

[bckr2549 / Breast-Cancer-Prediction](#) Public[Code](#) [Issues](#) [Pull requests](#) [Actions](#) [Projects](#) [Wiki](#) [Security](#) [Insights](#)[main](#) [...](#)[Breast-Cancer-Prediction / Breast Cancer Classification.ipynb](#)

bckr2549 Add files via upload

[History](#)

1 contributor

1583 lines (1583 sloc) | 62.1 KB

...

Importing the Dependencies

In [1]:

```
import numpy as np
import pandas as pd
import sklearn.datasets
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

Data Collection & Processing

In [2]:

```
# Loading the data from sklearn
breast_cancer_dataset = sklearn.datasets.load_breast_cancer()
```

In [3]:

```
print(breast_cancer_dataset)
```

```
1, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 0, 1, 0, 0, 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, 0, 0, 0, 0, 0, 0, 1]), 'frame': N
one, 'target_names': array(['malignant', 'benign'], dtype='
```

In [4]:

```
# Loading the data to a data frame
data_frame = pd.DataFrame(breast_cancer_dataset.data, columns = breast_cancer
```

In [5]:

```
# print the first 5 rows of the dataframe
data_frame.head()
```

Out[5]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	sym
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	

5 rows × 30 columns

In [6]:

```
# adding the 'target' column to the data frame
data_frame['label'] = breast_cancer_dataset.target
```

In [7]:

```
# print Last 5 rows of the dataframe
data_frame.tail()
```

Out[7]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	s)
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	

5 rows × 31 columns

In [8]:

```
# number of rows and columns in the dataset
data_frame.shape
```

Out[8]: (569, 31)

In [9]:

```
# getting some information about the data
data_frame.info()
```

RangeIndex: 569 entries, 0 to 568
Data columns (total 31 columns):

#	Column	Non-Null Count	Dtype
0	mean radius	569 non-null	float64
1	mean texture	569 non-null	float64
2	mean perimeter	569 non-null	float64
3	mean area	569 non-null	float64
4	mean smoothness	569 non-null	float64
5	mean compactness	569 non-null	float64
6	mean concavity	569 non-null	float64
7	mean concave points	569 non-null	float64
8	mean symmetry	569 non-null	float64
9	mean fractal dimension	569 non-null	float64
10	radius error	569 non-null	float64
11	texture error	569 non-null	float64
12	perimeter error	569 non-null	float64
13	area error	569 non-null	float64
14	smoothness error	569 non-null	float64
15	compactness error	569 non-null	float64
16	concavity error	569 non-null	float64
17	concave points error	569 non-null	float64
18	symmetry error	569 non-null	float64
19	fractal dimension error	569 non-null	float64
20	worst radius	569 non-null	float64
21	worst texture	569 non-null	float64
22	worst perimeter	569 non-null	float64
23	worst area	569 non-null	float64
24	worst smoothness	569 non-null	float64
25	worst compactness	569 non-null	float64
26	worst concavity	569 non-null	float64
27	worst concave points	569 non-null	float64
28	worst symmetry	569 non-null	float64
29	worst fractal dimension	569 non-null	float64
30	label	569 non-null	int32

dtypes: float64(30), int32(1)
memory usage: 135.7 KB

In [10]:

```
# statistical measures about the data
data_frame.describe()
```

Out[10]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	n conc
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.00
mean	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.081
std	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.071
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.001
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.021

50%	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.06
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.131
max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.421

8 rows × 31 columns

```
In [11]: # checking the distribution of Target Variable
data_frame['label'].value_counts()
```

```
Out[11]: 1    357
0    212
Name: label, dtype: int64
```

1 --> Benign

0 --> Malignant

```
In [12]: data_frame.groupby('label').mean()
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity
0	17.462830	21.604906	115.365377	978.376415	0.102898	0.145188	0.160775
1	12.146524	17.914762	78.075406	462.790196	0.092478	0.080085	0.046058

2 rows × 30 columns

Separating the features and target

```
In [13]: X = data_frame.drop(columns='label', axis=1)
Y = data_frame['label']
```

```
In [14]: print(X)
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness
0	17.99	10.38	122.80	1001.0	0.11840
1	20.57	17.77	132.90	1326.0	0.08474
2	19.69	21.25	130.00	1203.0	0.10960
3	11.42	20.38	77.58	386.1	0.14250
4	20.29	14.34	135.10	1297.0	0.10030
..
564	21.56	22.39	142.00	1479.0	0.11100

565	20.13	28.25	131.20	1261.0	0.09780
566	16.60	28.08	108.30	858.1	0.08455
567	20.60	29.33	140.10	1265.0	0.11780
568	7.76	24.54	47.92	181.0	0.05263
<hr/>					
0	mean compactness	mean concavity	mean concave points	mean symmetry	\
1	0.27760	0.30010	0.14710	0.2419	
2	0.07864	0.08690	0.07017	0.1812	
3	0.15990	0.19740	0.12790	0.2069	
4	0.28390	0.24140	0.10520	0.2597	
..	0.13280	0.19800	0.10430	0.1809	
564	0.11590	0.24390	0.13890	0.1726	
565	0.10340	0.14400	0.09791	0.1752	
566	0.10230	0.09251	0.05302	0.1590	
567	0.27700	0.35140	0.15200	0.2397	
568	0.04362	0.00000	0.00000	0.1587	
<hr/>					
0	mean fractal dimension	...	worst radius	worst texture	\
1	0.07871	...	25.380	17.33	
2	0.05667	...	24.990	23.41	
3	0.05999	...	23.570	25.53	
4	0.09744	...	14.910	26.50	
..	0.05883	...	22.540	16.67	
564	0.05623	...	25.450	26.40	
565	0.05533	...	23.690	38.25	
566	0.05648	...	18.980	34.12	
567	0.07016	...	25.740	39.42	
568	0.05884	...	9.456	30.37	
<hr/>					
0	worst perimeter	worst area	worst smoothness	worst compactness	\
1	184.60	2019.0	0.16220	0.66560	
2	158.80	1956.0	0.12380	0.18660	
3	152.50	1709.0	0.14440	0.42450	
4	98.87	567.7	0.20980	0.86630	
..	152.20	1575.0	0.13740	0.20500	
564	166.10	2027.0	0.14100	0.21130	
565	155.00	1731.0	0.11660	0.19220	
566	126.70	1124.0	0.11390	0.30940	
567	184.60	1821.0	0.16500	0.86810	
568	59.16	268.6	0.08996	0.06444	
<hr/>					
0	worst concavity	worst concave points	worst symmetry	\	
1	0.7119	0.2654	0.4601		
2	0.2416	0.1860	0.2750		
3	0.4504	0.2430	0.3613		
4	0.6869	0.2575	0.6638		
..	0.4000	0.1625	0.2364		
564	0.4107	0.2216	0.2060		
565	0.3215	0.1628	0.2572		
566	0.3403	0.1418	0.2218		
567	0.9387	0.2650	0.4087		
568	0.0000	0.0000	0.2871		
<hr/>					
0	worst fractal dimension				
1	0.11890				
1	0.08902				

```
2          0.08758
3          0.17300
4          0.07678
..
564         ...
565         0.06637
566         0.07820
567         0.12400
568         0.07039
```

[569 rows x 30 columns]

```
In [15]: print(Y)
```

```
0      0
1      0
2      0
3      0
4      0
..
564    0
565    0
566    0
567    0
568    1
Name: label, Length: 569, dtype: int32
```

Splitting the data into training data & Testing data

```
In [16]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, ran
```

```
In [17]: print(X.shape, X_train.shape, X_test.shape)
```

(569, 30) (455, 30) (114, 30)

Model Training

Logistic Regression

```
In [18]: model = LogisticRegression()
```

```
In [19]: # training the Logistic Regression model using Training data
model.fit(X_train, Y_train)
```

```
C:\Users\sathy\anaconda3\lib\site-packages\sklearn\linear_model\_logistic.p
y:814: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.
```

```
Increase the number of iterations (max_iter) or scale the data as shown in:  
    https://scikit-learn.org/stable/modules/preprocessing.html  
Please also refer to the documentation for alternative solver options:  
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regre  
ssion  
    n_iter_i = _check_optimize_result()  
Out[19]: LogisticRegression()
```

Model Evaluation

Accuracy Score

```
In [20]:  
# accuracy on training data  
X_train_prediction = model.predict(X_train)  
training_data_accuracy = accuracy_score(Y_train, X_train_prediction)
```

```
In [21]:  
print('Accuracy on training data = ', training_data_accuracy)
```

```
Accuracy on training data =  0.9384615384615385
```

```
In [22]:  
# accuracy on test data  
X_test_prediction = model.predict(X_test)  
test_data_accuracy = accuracy_score(Y_test, X_test_prediction)
```

```
In [23]:  
print('Accuracy on test data = ', test_data_accuracy)
```

```
Accuracy on test data =  0.9298245614035088
```

Building a Predictive System

```
In [24]:  
input_data = (13.54, 14.36, 87.46, 566.3, 0.09779, 0.08129, 0.06664, 0.04781, 0.1885  
  
# change the input data to a numpy array  
input_data_as_numpy_array = np.asarray(input_data)  
  
# reshape the numpy array as we are predicting for one datapoint  
input_data_reshaped = input_data_as_numpy_array.reshape(1,-1)  
  
prediction = model.predict(input_data_reshaped)  
print(prediction)  
  
if (prediction[0] == 0):  
    print('The Breast cancer is Malignant')  
  
else:  
    print('The Breast Cancer is Benign')
```