

[illegible]

```

0.112 2.873\n texture (standard error): 0.36 4.885\n perimeter (standard error):
0.757 21.98\n area (standard error): 6.802 542.2\n smoothness (standard error):
0.002 0.031\n compactness (standard error): 0.002 0.135\n concavity (standard error):
0.0 0.396\n concave points (standard error): 0.0 0.053\n symmetry (standard error):
0.008 0.079\n fractal dimension (standard error): 0.001 0.03\n radius (worst):
7.93 36.04\n texture (worst): 12.02 49.54\n perimeter (worst):
50.41 251.2\n area (worst): 185.2 4254.0\n smoothness (worst):
0.071 0.223\n compactness (worst): 0.027 1.058\n concavity (worst):
0.0 1.252\n concave points (worst): 0.0 0.291\n symmetry (worst):
0.156 0.664\n fractal dimension (worst): 0.055 0.208\n =====
===== \n\n :Missing Attribute Values: None\n\n :Class Distribution: 212 - Malignant, 357 - Benign\n
\n\n :Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian\n\n :Donor: Nick Street\n\n
:Date: November, 1995\n\nThis is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.\nhttps://goo.g
l/U2Uwz2\n\nFeatures are computed from a digitized image of a fine needle\naspirate (FNA) of a breast mass. Th
ey describe\ncharacteristics of the cell nuclei present in the image.\n\nSeparating plane described above was o
btained using\nMultisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree\nConstruction Via Linear Program
ming." Proceedings of the 4th\nMidwest Artificial Intelligence and Cognitive Science Society,\npp. 97-101, 1992
], a classification method which uses linear\nprogramming to construct a decision tree. Relevant features\nwer
e selected using an exhaustive search in the space of 1-4\nfeatures and 1-3 separating planes.\n\nThe actual li
near program used to obtain the separating plane\nin the 3-dimensional space is that described in:\n[K. P. Benn
ett and O. L. Mangasarian: "Robust Linear\nProgramming Discrimination of Two Linearly Inseparable Sets",\nOptim
ization Methods and Software 1, 1992, 23-34].\n\nThis database is also available through the UW CS ftp server:\
\n\nftp ftp.cs.wisc.edu\nncd math-prog/cpo-dataset/machine-learn/WDBC/\n\n.. topic:: References\n\n - W.N. Stre
et, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction \n for breast tumor diagnosis. IS&T/SPIE
1993 International Symposium on \n Electronic Imaging: Science and Technology, volume 1905, pages 861-870,\
\n San Jose, CA, 1993.\n - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and \n
prognosis via linear programming. Operations Research, 43(4), pages 570-577, \n July-August 1995.\n - W.H
. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques\n to diagnose breast cancer from
fine-needle aspirates. Cancer Letters 77 (1994) \n 163-171.', 'feature_names': array(['mean radius', 'mean
texture', 'mean perimeter', 'mean area',
'mean smoothness', 'mean compactness', 'mean concavity',
'mean concave points', 'mean symmetry', 'mean fractal dimension',
'radius error', 'texture error', 'perimeter error', 'area error',
'smoothness error', 'compactness error', 'concavity error',
'concave points error', 'symmetry error',
'fractal dimension error', 'worst radius', 'worst texture',
'worst perimeter', 'worst area', 'worst smoothness',
'worst compactness', 'worst concavity', 'worst concave points',
'worst symmetry', 'worst fractal dimension'], dtype='<U23'), 'filename': 'breast_cancer.csv', 'data_modu
le': 'sklearn.datasets.data'}

```

```
In [12]: data_frame=pd.DataFrame(breast_cancer_dataset.data,columns=breast_cancer_dataset.feature_names)
```

```
In [13]: data_frame.head()
```

```
Out[13]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius	worst texture	worst perimeter	
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	25.38	17.33	184.60	2
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	24.99	23.41	158.80	1
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	23.57	25.53	152.50	1
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	14.91	26.50	98.87	
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	22.54	16.67	152.20	1

5 rows × 30 columns

## Adding the 'Target' column to the data frame

```
In [14]: data_frame['Label']=breast_cancer_dataset.target
```

```
In [15]: data_frame.tail()
```

```
Out[15]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst texture	worst perimeter	worst area
564	21.56	22.39	142.00	1479.0	0.11100	0.11590	0.24390	0.13890	0.1726	0.05623	...	26.40	166.10	2027.0
565	20.13	28.25	131.20	1261.0	0.09780	0.10340	0.14400	0.09791	0.1752	0.05533	...	38.25	155.00	1731.0
566	16.60	28.08	108.30	858.1	0.08455	0.10230	0.09251	0.05302	0.1590	0.05648	...	34.12	126.70	1124.0
567	20.60	29.33	140.10	1265.0	0.11780	0.27700	0.35140	0.15200	0.2397	0.07016	...	39.42	184.60	1821.0
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.00000	0.00000	0.1587	0.05884	...	30.37	59.16	268.6

5 rows × 31 columns

0 ----> Malignant 1 ----> Benign

# Number of rows and columns in data frame

```
In [16]: data_frame.shape
```

```
Out[16]: (569, 31)
```

## Checking for null values in the data frame

```
In [17]: data_frame.isnull().sum()
```

```
Out[17]: mean radius          0
mean texture          0
mean perimeter        0
mean area             0
mean smoothness       0
mean compactness      0
mean concavity         0
mean concave points   0
mean symmetry         0
mean fractal dimension 0
radius error          0
texture error         0
perimeter error       0
area error            0
smoothness error      0
compactness error     0
concavity error       0
concave points error  0
symmetry error        0
fractal dimension error 0
worst radius          0
worst texture         0
worst perimeter       0
worst area            0
worst smoothness      0
worst compactness     0
worst concavity       0
worst concave points  0
worst symmetry        0
worst fractal dimension 0
Label                0
dtype: int64
```

There are no null values in the dataframe.

## Statistical insights of the data frame

```
In [19]: data_frame.describe()
```

```
Out[19]:
```

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	
<b>count</b>	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.
<b>mean</b>	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.088799	0.048919	0.181162	0.062798	...	25.
<b>std</b>	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079720	0.038803	0.027414	0.007060	...	6.
<b>min</b>	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000	0.106000	0.049960	...	12.
<b>25%</b>	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310	0.161900	0.057700	...	21.
<b>50%</b>	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061540	0.033500	0.179200	0.061540	...	25.
<b>75%</b>	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074000	0.195700	0.066120	...	29.
<b>max</b>	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200	0.304000	0.097440	...	49.

8 rows × 31 columns

## Checking the distribution of Target Variable

```
In [20]: data_frame['Label'].value_counts()
```

```
Out[20]: 1    357
0     212
Name: Label, dtype: int64
```

In the given dataset there are 357 benign cases and 212 malignant cases of breast cancer.

# Grouping the data frame according to 'Label'

```
In [22]: data_frame.groupby('Label').mean()
```

Out[22]:

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points	mean symmetry	mean fractal dimension	...	worst radius
Label												
0	17.462830	21.604906	115.365377	978.376415	0.102898	0.145188	0.160775	0.087990	0.192909	0.062680	...	21.134811
1	12.146524	17.914762	78.075406	462.790196	0.092478	0.080085	0.046058	0.025717	0.174186	0.062867	...	13.379801

2 rows × 30 columns

In malignant cases almost all parameters are of greater value than those of benign cases.

## Seperating features and 'Label'

```
In [23]: X=data_frame.drop(columns='Label',axis=1)
Y=data_frame['Label']
```

```
In [24]: 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

	mean compactness	mean concavity	mean concave points	mean symmetry \
0	0.27760	0.30010	0.14710	0.2419
1	0.07864	0.08690	0.07017	0.1812
2	0.15990	0.19740	0.12790	0.2069
3	0.28390	0.24140	0.10520	0.2597
4	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

	mean fractal dimension	...	worst radius	worst texture \
0	0.07871	...	25.380	17.33
1	0.05667	...	24.990	23.41
2	0.05999	...	23.570	25.53
3	0.09744	...	14.910	26.50
4	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

	worst perimeter	worst area	worst smoothness	worst compactness \
0	184.60	2019.0	0.16220	0.66560
1	158.80	1956.0	0.12380	0.18660
2	152.50	1709.0	0.14440	0.42450
3	98.87	567.7	0.20980	0.86630
4	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

	worst concavity	worst concave points	worst symmetry \
0	0.7119	0.2654	0.4601
1	0.2416	0.1860	0.2750
2	0.4504	0.2430	0.3613
3	0.6869	0.2575	0.6638
4	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

	worst fractal dimension
0	0.11890
1	0.08902
2	0.08758
3	0.17300
4	0.07678
..	...
564	0.07115
565	0.06637
566	0.07820
567	0.12400
568	0.07039

[569 rows x 30 columns]

In [25]: 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 & Test Data

```
In [26]: X_train , X_test , Y_train , Y_test =train_test_split(X,Y,test_size=0.2,random_state=2)
```

```
In [27]: print(X.shape,X_train.shape,X_test.shape)

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

## Model Training using Logistic Regression

```
In [30]: model=LogisticRegression()
```

```
In [31]: model.fit(X_train,Y_train)
```

```

C:\ProgramData\anaconda3\Lib\site-packages\sklearn\linear_model\_logistic.py:460: 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-regression
n_iter_i = _check_optimize_result(

```

```
Out[31]: LogisticRegression
LogisticRegression()
```

## Model Evaluation

### Accuracy Score

### Accuracy on Training Data

```
In [35]: X_train_prediction=model.predict(X_train)
training_data_accuracy=accuracy_score(Y_train,X_train_prediction)
```

```
In [36]: print('Accuracy on Training Data=',training_data_accuracy)

Accuracy on Training Data= 0.9494505494505494
```

### Accuracy on Test Data

```
In [37]: X_test_prediction=model.predict(X_test)
test_data_accuracy=accuracy_score(Y_test,X_test_prediction)
print('Accuracy on Test Data=',test_data_accuracy)

Accuracy on Test Data= 0.9210526315789473
```

## Building a Predictive System

```
In [38]: input_data=(17.99,10.38,122.8,1001,0.1184,0.2776,0.3001,0.1471,0.2419,0.07871,1.095,0.9053,8.589,153.4,0.006399)
```

## Changing the input data into a Numpy array

```
In [39]: input_data_numpy_format=np.asarray(input_data)
```

# Reshaping the Numpy Array as we are predicting for one datapoint

```
In [40]: input_data_reshape=input_data_numpy_format.reshape(1,-1)
```

## The Prediction

```
In [42]: prediction=model.predict(input_data_reshape)
print(prediction)
```

```
[0]
```

```
C:\ProgramData\anaconda3\Lib\site-packages\sklearn\base.py:464: UserWarning: X does not have valid feature names, but LogisticRegression was fitted with feature names
warnings.warn(
```

```
In [44]: if (prediction[0] == 0):
print("THIS IS A CASE OF MALIGNANT BREAST CANCER.")
else:
print("This IS A CASE OF BENIGN BREAST CANCER.")
```

```
THIS IS A CASE OF MALIGNANT BREAST CANCER.
```

```
In [ ]:
```

```
Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js
```