

```
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
```

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

[illegible]

```
# loading the data to a data frame
data_frame = pd.DataFrame(breast_cancer_dataset.data, columns = breast_cancer_dataset.feature_names)
# print the first 5 rows of the dataframe
data_frame.head()
```

1 to 5 of 5 entries

Filter

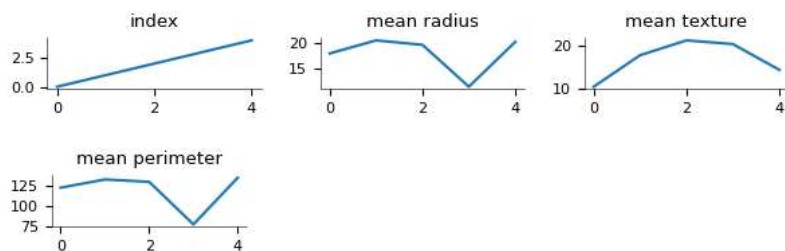
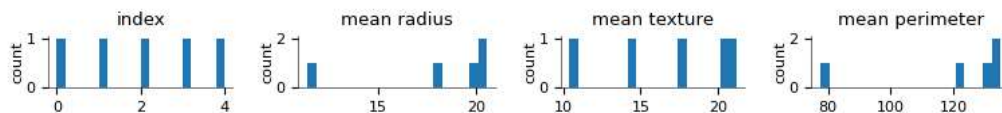


index	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity
0	17.99	10.38	122.8	1001.0	0.1184	0.2776	0.3001
1	20.57	17.77	132.9	1326.0	0.08474	0.07864	0.0869
2	19.69	21.25	130.0	1203.0	0.1096	0.1599	0.1974
3	11.42	20.38	77.58	386.1	0.1425	0.2839	0.2414
4	20.29	14.34	135.1	1297.0	0.1003	0.1328	0.198

Show 25 per page

Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

Warning: Total number of columns (30) exceeds max_columns (20) limiting to first (20) columns.

Values**Distributions****2-d distributions**

```
# adding the 'target' column to the data frame
data_frame['label'] = breast_cancer_dataset.target
# print last 5 rows of the dataframe
data_frame.tail()
```

1 to 5 of 5 entries

Filter

?

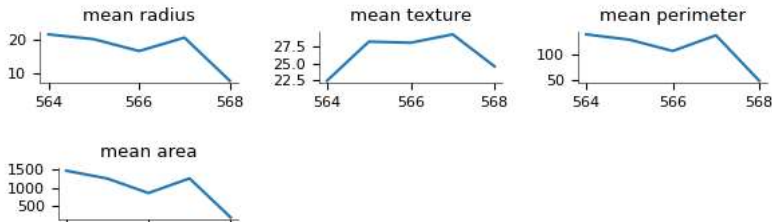
index	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity
564	21.56	22.39	142.0	1479.0	0.111	0.1159	0.2439
565	20.13	28.25	131.2	1261.0	0.0978	0.1034	0.144
566	16.6	28.08	108.3	858.1	0.08455	0.1023	0.09251
567	20.6	29.33	140.1	1265.0	0.1178	0.277	0.3514
568	7.76	24.54	47.92	181.0	0.05263	0.04362	0.0

Show 25 per page



Like what you see? Visit the [data table notebook](#) to learn more about interactive tables.

Values



```
# number of rows and columns in the dataset
data_frame.shape
# getting some information about the data
data_frame.info()
```

```
<class 'pandas.core.frame.DataFrame'>
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    int64
dtypes: float64(30), int64(1)
memory usage: 137.9 KB
```

```
# checking for missing values
data_frame.isnull().sum()

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
```

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

	mean radius	mean texture	mean perimeter	mean area	mean smoothness	mean compactness	mean concavity	mean concave points
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000
mean	14.127292	19.289649	91.969033	654.889104	0.096360	0.104341	0.088799	0.048919
std	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079720	0.038803
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310
50%	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061540	0.033500
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074000
max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200

8 rows × 31 columns



```
# checking the distribution of Target Varibale
data_frame['label'].value_counts()

1    357
0    212
Name: label, dtype: int64
```

1 --> Benign

0 --> Malignant

```
data_frame.groupby('label').mean()
```

Separating the features and target

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

```
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]
```

```
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: int64
```

Splitting the data into training data & Testing data

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=2)
print(X.shape, X_train.shape, X_test.shape)
```

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

Model Training**Logistic Regression**

```
model = LogisticRegression()
```

```

model = LogisticRegression()
# training the Logistic Regression model using Training data
model.fit(X_train, Y_train)

/usr/local/lib/python3.10/dist-packages/sklearn/linear_model/_logistic.py:458: ConvergenceWarning: lbfgs
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(
  ▾ LogisticRegression
  LogisticRegression()

```

Model Evaluation

Accuracy Score

```

# accuracy on training data
X_train_prediction = model.predict(X_train)
training_data_accuracy = accuracy_score(Y_train, X_train_prediction)
print('Accuracy on training data = ', training_data_accuracy)
# accuracy on test data
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 training data = 0.9472527472527472
Accuracy on test data = 0.9298245614035088

```

Building a Predictive System

```

input_data = (13.54,14.36,87.46,566.3,0.09779,0.08129,0.06664,0.04781,0.1885,0.05766,0.2699,0.7886,2.058,23.56,0.008462,0.0146,0.02387,0.0131)

# 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')

[1]
The Breast Cancer is Benign
/usr/local/lib/python3.10/dist-packages/sklearn/base.py:439: UserWarning: X does not have valid feature names, but LogisticRegression will
warnings.warn(

```

✓ 0s completed at 12:22 PM

● ✕