

Breast Cancer Detection Machine Learning End to End Project

Import essential libraries

```
1 | # import libraries
2 | import pandas as pd # for data manipulation or analysis
3 | import numpy as np # for numeric calculation
4 | import matplotlib.pyplot as plt # for data visualization
5 | import seaborn as sns # for data visualization
```

Load breast cancer dataset & explore

```
1 | #Load breast cancer dataset
2 | from sklearn.datasets import load_breast_cancer
3 | cancer_dataset = load_breast_cancer()
```

```
1 | type(cancer_dataset)
```

Output >>> sklearn.utils.Bunch
The scikit-learn store data in an object bunch like a dictionary.

```
1 | # keys in dataset
2 | cancer_dataset.keys()
```

Output >>> dict_keys(['data', 'target', 'target_names', 'DESCR', 'feature_names', 'filename'])

```
1 | # featur of each cells in numeric format
2 | cancer_dataset['data']
```

Output >>>

```
1 | array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01,
2 |        1.189e-01],
3 |        [2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01,
4 |        8.902e-02],
5 |        [1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01,
6 |        8.758e-02],
7 |        ...,
8 |        [1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,
9 |        7.820e-02],
10 |        [2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01,
11 |        1.240e-01],
12 |        [7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01,
13 |        7.039e-02]])
```

These numeric values are extracted features of each cell.

```
1 | # malignant or benign value
2 | cancer_dataset['target']
```

The **target** stores the values of malignant or benign tumors.

```
1 | # target value name malignant or benign tumor
2 | cancer_dataset['target_names']
```

Output >>> array(['malignant', 'benign'], dtype='<U9')
0 means malignant tumor
1 mean benign tumor

The **cancer_dataset['DESCR']** store the description of breast cancer dataset.

```
1 | # description of data
2 | print(cancer_dataset['DESCR'])
```

Output >>>

```
1 | .. _breast_cancer_dataset:
2 |
3 | Breast cancer wisconsin (diagnostic) dataset
```

```
4 -----
5
6 **Data Set Characteristics:**
7   :Number of Instances: 569
8   :Number of Attributes: 30 numeric, predictive attributes and the class
9   :Attribute Information:
10      - radius (mean of distances from center to points on the perimeter)
11      - texture (standard deviation of gray-scale values)
12      - perimeter
13      - area
14      - smoothness (local variation in radius lengths)
15      - compactness (perimeter^2 / area - 1.0)
16      - concavity (severity of concave portions of the contour)
17      - concave points (number of concave portions of the contour)
18      - symmetry
19      - fractal dimension ("coastline approximation" - 1)
20
21   The mean, standard error, and "worst" or largest (mean of the three
22   largest values) of these features were computed for each image,
23   resulting in 30 features.  For instance, field 3 is Mean Radius, field
24   13 is Radius SE, field 23 is Worst Radius.
25
26   - class:
27     - WDBC-Malignant
28     - WDBC-Benign
29
30   :Summary Statistics:
31   =====
32                                     Min      Max
33   =====
34   radius (mean):                    6.981   28.11
35   texture (mean):                   9.71    39.28
36   perimeter (mean):                 43.79   188.5
37   area (mean):                      143.5   2501.0
38   smoothness (mean):                0.053   0.163
39   compactness (mean):               0.019   0.345
40   concavity (mean):                 0.0     0.427
41   concave points (mean):            0.0     0.201
42   symmetry (mean):                  0.106   0.304
43   fractal dimension (mean):         0.05    0.097
44   radius (standard error):          0.112   2.873
45   texture (standard error):         0.36    4.885
46   perimeter (standard error):       0.757   21.98
47   area (standard error):            6.802   542.2
48   smoothness (standard error):      0.002   0.031
49   compactness (standard error):     0.002   0.135
50   concavity (standard error):       0.0     0.396
51   concave points (standard error):  0.0     0.053
52   symmetry (standard error):        0.008   0.079
53   fractal dimension (standard error): 0.001   0.03
54   radius (worst):                   7.93    36.04
55   texture (worst):                  12.02   49.54
56   perimeter (worst):                50.41   251.2
57   area (worst):                     185.2   4254.0
58   smoothness (worst):               0.071   0.223
59   compactness (worst):               0.027   1.058
60   concavity (worst):                0.0     1.252
61   concave points (worst):           0.0     0.291
62   symmetry (worst):                 0.156   0.664
63   fractal dimension (worst):        0.055   0.208
64   =====
65
66   :Missing Attribute Values: None
67   :Class Distribution: 212 - Malignant, 357 - Benign
68   :Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian
69   :Donor: Nick Street
70   :Date: November, 1995
71
72   This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.
73   https://goo.gl/U2Uwz2
74
75   Features are computed from a digitized image of a fine needle
76   aspirate (FNA) of a breast mass. They describe
77   characteristics of the cell nuclei present in the image.
78
79   Separating plane described above was obtained using
80   Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree
81   Construction Via Linear Programming." Proceedings of the 4th
82   Midwest Artificial Intelligence and Cognitive Science Society,
83   pp. 97-101, 1992], a classification method which uses linear
84   programming to construct a decision tree. Relevant features
85   were selected using an exhaustive search in the space of 1-4
86   features and 1-3 separating planes.
87
88   The actual linear program used to obtain the separating plane
89   in the 3-dimensional space is that described in:
90   [K. P. Bennett and O. L. Mangasarian: "Robust Linear
91   Programming Discrimination of Two Linearly Inseparable Sets",
92   Optimization Methods and Software 1, 1992, 23-34].
93
94   This database is also available through the UW CS ftp server:
95
96   ftp ftp.cs.wisc.edu
97   cd math-prog/cpo-dataset/machine-learn/WDBC/
98
99   .. topic:: References
100
101   - W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction
102     for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on
103     Electronic Imaging: Science and Technology, volume 1905, pages 861-870,
104     San Jose, CA, 1993.
105   - O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and
106     prognosis via linear programming. Operations Research, 43(4), pages 570-577,
107     July-August 1995.
108   - W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques
109     to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994)
```

Features name of malignant & benign tumor.

```
1 | # name of features
2 | print(cancer_dataset['feature_names'])
```

Output >>>

```
1 | ['mean radius' 'mean texture' 'mean perimeter' 'mean area'
2 | 'mean smoothness' 'mean compactness' 'mean concavity'
3 | 'mean concave points' 'mean symmetry' 'mean fractal dimension'
4 | 'radius error' 'texture error' 'perimeter error' 'area error'
5 | 'smoothness error' 'compactness error' 'concavity error'
6 | 'concave points error' 'symmetry error' 'fractal dimension error'
7 | 'worst radius' 'worst texture' 'worst perimeter' 'worst area'
8 | 'worst smoothness' 'worst compactness' 'worst concavity'
9 | 'worst concave points' 'worst symmetry' 'worst fractal dimension']
```

```
1 | # location/path of data file
2 | print(cancer_dataset['filename'])
```

Output >>> C:\ProgramData\Anaconda3\lib\site-packages\sklearn\datasets\data\breast_cancer.csv

Create DataFrame

```
1 | # create datafrmae
2 | cancer_df = pd.DataFrame(np.c_[cancer_dataset['data'],cancer_dataset['target']],
3 |                          columns = np.append(cancer_dataset['feature_names'], ['target']))
```

Head of cancer DataFrame

```
1 | # Head of cancer DataFrame
2 | cancer_df.head(6)
```

Output >>>

	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	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension	target
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	0.2419	0.07871	...	17.33	184.60	2019.0	0.1622	0.6656	0.7119	0.2654	0.4601	0.11890	0.0
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	0.1812	0.05667	...	23.41	158.80	1956.0	0.1238	0.1866	0.2416	0.1860	0.2750	0.08902	0.0
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	0.2069	0.05999	...	25.53	152.50	1709.0	0.1444	0.4245	0.4504	0.2430	0.3613	0.08758	0.0
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	0.2597	0.09744	...	26.50	98.87	567.7	0.2098	0.8663	0.6869	0.2575	0.6638	0.17300	0.0
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	0.1809	0.05883	...	16.67	152.20	1575.0	0.1374	0.2050	0.4000	0.1625	0.2364	0.07678	0.0
5	12.45	15.70	82.57	477.1	0.12780	0.17000	0.1578	0.08089	0.2087	0.07613	...	23.75	103.40	741.6	0.1791	0.5249	0.5355	0.1741	0.3985	0.12440	0.0
6 rows x 31 columns																					

The tail of cancer DataFrame

```
1 | # Tail of cancer DataFrame
2 | cancer_df.tail(6)
```

Output >>>

	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	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension	target
563	20.92	25.09	143.00	1347.0	0.10990	0.22360	0.31740	0.14740	0.2149	0.06879	...	29.41	179.10	1819.0	0.14070	0.41860	0.6599	0.2542	0.2929	0.09873	0.0
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	0.14100	0.21130	0.4107	0.2216	0.2060	0.07115	0.0
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	0.11660	0.19220	0.3215	0.1628	0.2572	0.06637	0.0
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	0.11390	0.30940	0.3403	0.1418	0.2218	0.07820	0.0
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	0.16500	0.86810	0.9387	0.2650	0.4087	0.12400	0.0
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	0.08996	0.06444	0.0000	0.0000	0.2871	0.07039	1.0
6 rows x 31 columns																					

Getting information of cancer DataFrame using '.info()' method.

```
1 | # Information of cancer Dataframe
2 | cancer_df.info()
```

Output >>>

```
1 | <class 'pandas.core.frame.DataFrame'>
2 | RangeIndex: 569 entries, 0 to 568
3 | Data columns (total 31 columns):
4 | mean radius                569 non-null float64
5 | mean texture                569 non-null float64
6 | mean perimeter              569 non-null float64
```

```
7 | mean area          569 non-null float64
8 | mean smoothness    569 non-null float64
9 | mean compactness    569 non-null float64
10 | mean concavity       569 non-null float64
11 | mean concave points  569 non-null float64
12 | mean symmetry        569 non-null float64
13 | mean fractal dimension 569 non-null float64
14 | radius error         569 non-null float64
15 | texture error        569 non-null float64
16 | perimeter error      569 non-null float64
17 | area error           569 non-null float64
18 | smoothness error     569 non-null float64
19 | compactness error    569 non-null float64
20 | concavity error      569 non-null float64
21 | concave points error  569 non-null float64
22 | symmetry error       569 non-null float64
23 | fractal dimension error 569 non-null float64
24 | worst radius         569 non-null float64
25 | worst texture         569 non-null float64
26 | worst perimeter      569 non-null float64
27 | worst area           569 non-null float64
28 | worst smoothness     569 non-null float64
29 | worst compactness    569 non-null float64
30 | worst concavity       569 non-null float64
31 | worst concave points  569 non-null float64
32 | worst symmetry        569 non-null float64
33 | worst fractal dimension 569 non-null float64
34 | target               569 non-null float64
35 | dtypes: float64(31)
36 | memory usage: 137.9 KB
```

Total of non-null 569 patients’ information with 31 features. All feature data types in the float. The size of the DataFrame is 137.9 KB.

Numerical distribution of data. We can know to mean, standard deviation, min, max, 25%,50% and 75% value of each feature.

```
1 | # Numerical distribution of data
2 | cancer_df.describe()
```

Output >>>

	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	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension	target
count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.000000	569.000000	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	0.181162	0.062798	...	25.677223	107.261213	880.583128	0.132369	0.254265	0.272188	0.114606	0.290076	0.083946	0.627417
std	3.524049	4.301036	24.298981	351.914129	0.014064	0.052813	0.079720	0.038803	0.027414	0.007060	...	6.146258	33.602542	569.356993	0.022832	0.157336	0.208624	0.065732	0.061867	0.018061	0.483918
min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000	0.106000	0.049960	...	12.020000	50.410000	185.200000	0.071170	0.027290	0.000000	0.000000	0.156500	0.055040	0.000000
25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064920	0.029560	0.020310	0.161900	0.057700	...	21.080000	84.110000	515.300000	0.116600	0.147200	0.114500	0.064930	0.250400	0.071460	0.000000
50%	13.370000	18.840000	86.240000	551.100000	0.095870	0.092630	0.061540	0.033500	0.179200	0.061540	...	25.410000	97.660000	686.500000	0.131300	0.211900	0.226700	0.099930	0.282200	0.080040	1.000000
75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.130400	0.130700	0.074000	0.195700	0.066120	...	29.720000	125.400000	1084.000000	0.146000	0.339100	0.382900	0.161400	0.317900	0.092080	1.000000
max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200	0.304000	0.097440	...	49.540000	251.200000	4254.000000	0.222600	1.058000	1.252000	0.291000	0.663800	0.207500	1.000000

8 rows x 31 columns

Data Visualization

Pair plot of breast cancer data

Basically, the pair plot is used to show the numeric distribution in the scatter plot.

```
1 | # Paiplot of cancer dataframe
2 | sns.pairplot(cancer_df, hue = 'target')
```

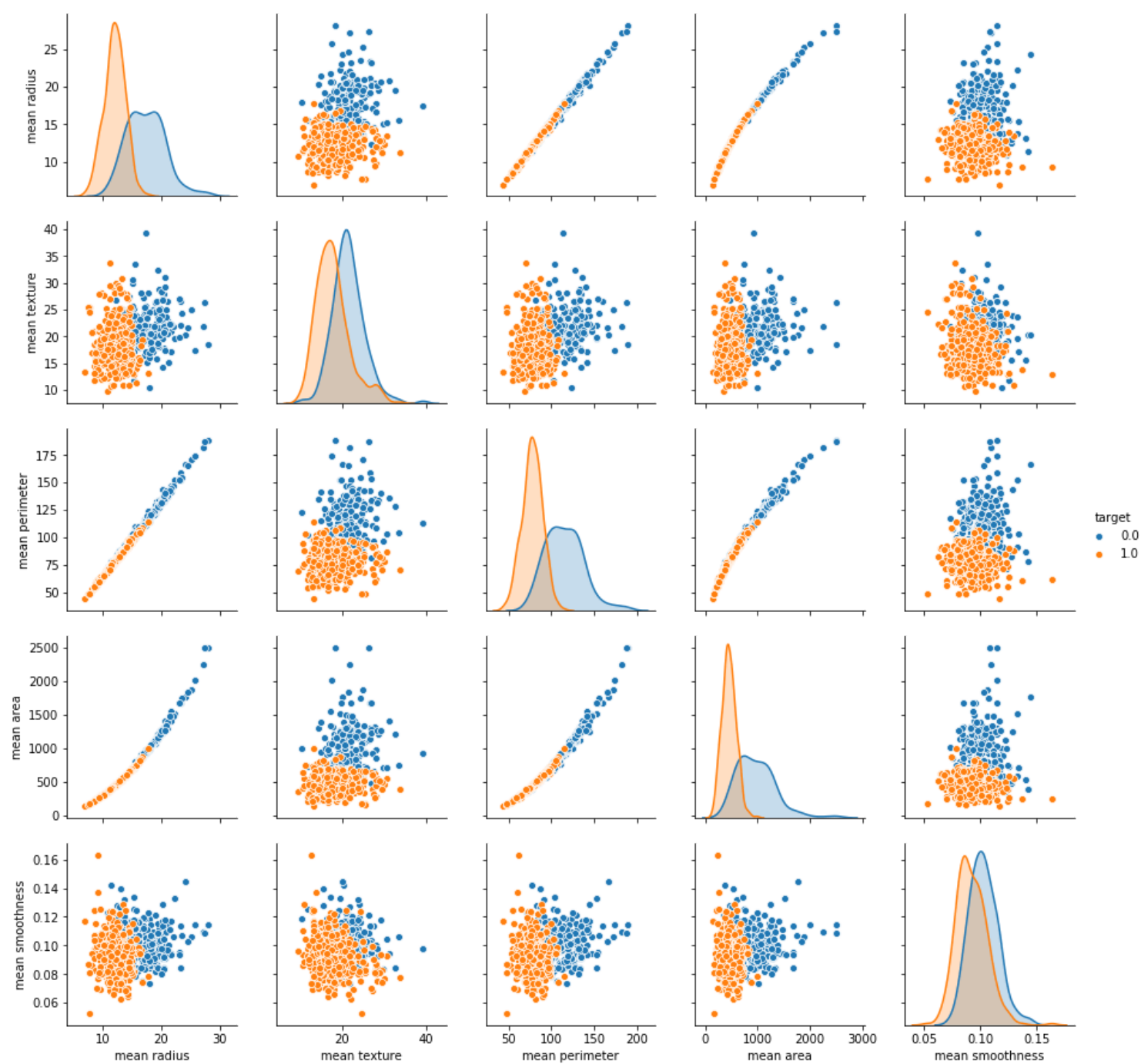
Output >>>



Pair plot of sample feature of DataFrame

```
1 # pair plot of sample feature
2 sns.pairplot(cancer_df, hue = 'target',
3             vars = ['mean radius', 'mean texture', 'mean perimeter', 'mean area', 'mean smoothness'] )
```

Output >>>



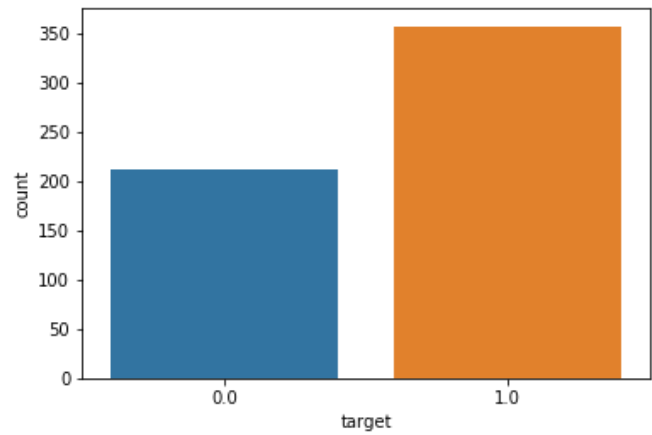
The pair plot showing malignant and benign tumor data distributed in two classes. It is easy to differentiate in the pair plot.

Counterplot

Showing the total count of malignant and benign tumor patients in counterplot.

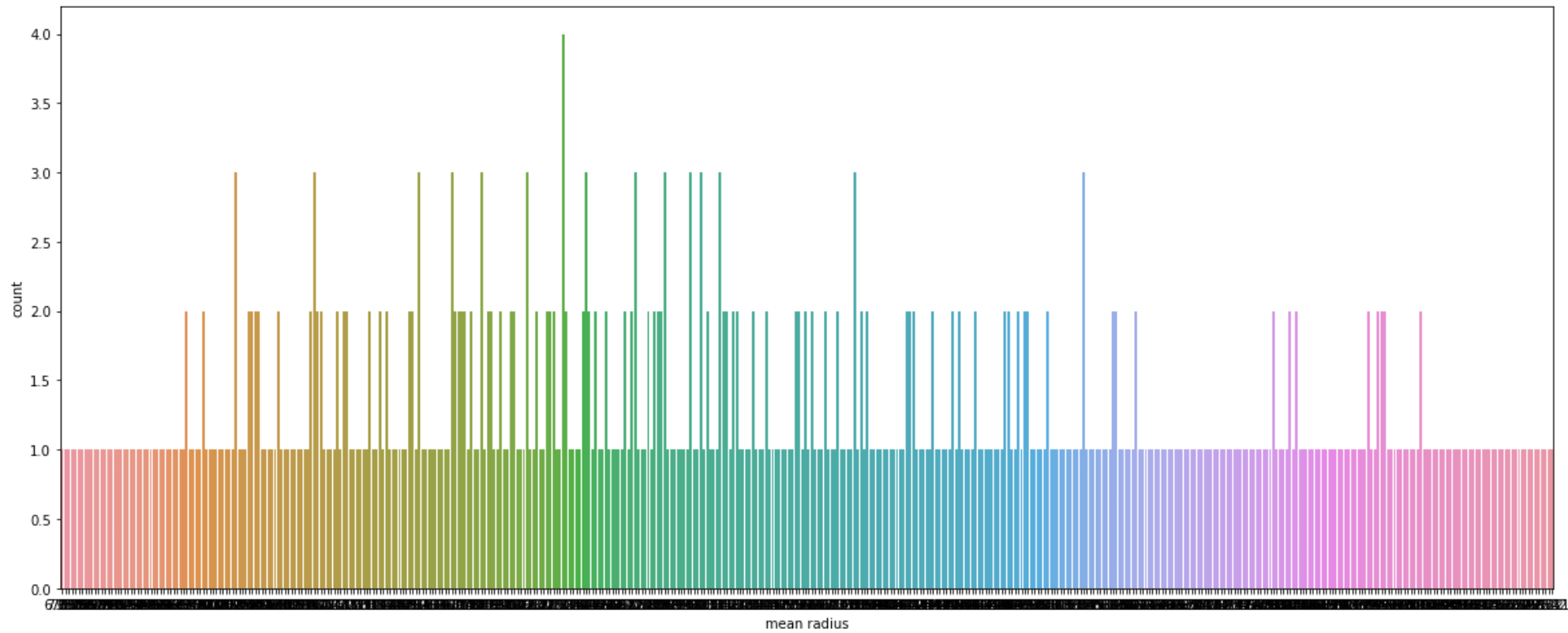
```
1 # Count the target class
2 sns.countplot(cancer_df['target'])
```

Output >>>



In the below counterplot max samples mean radius is equal to 1.

```
1 # counter plot of feature mean radius
2 plt.figure(figsize = (20,8))
3 sns.countplot(cancer_df['mean radius'])
```



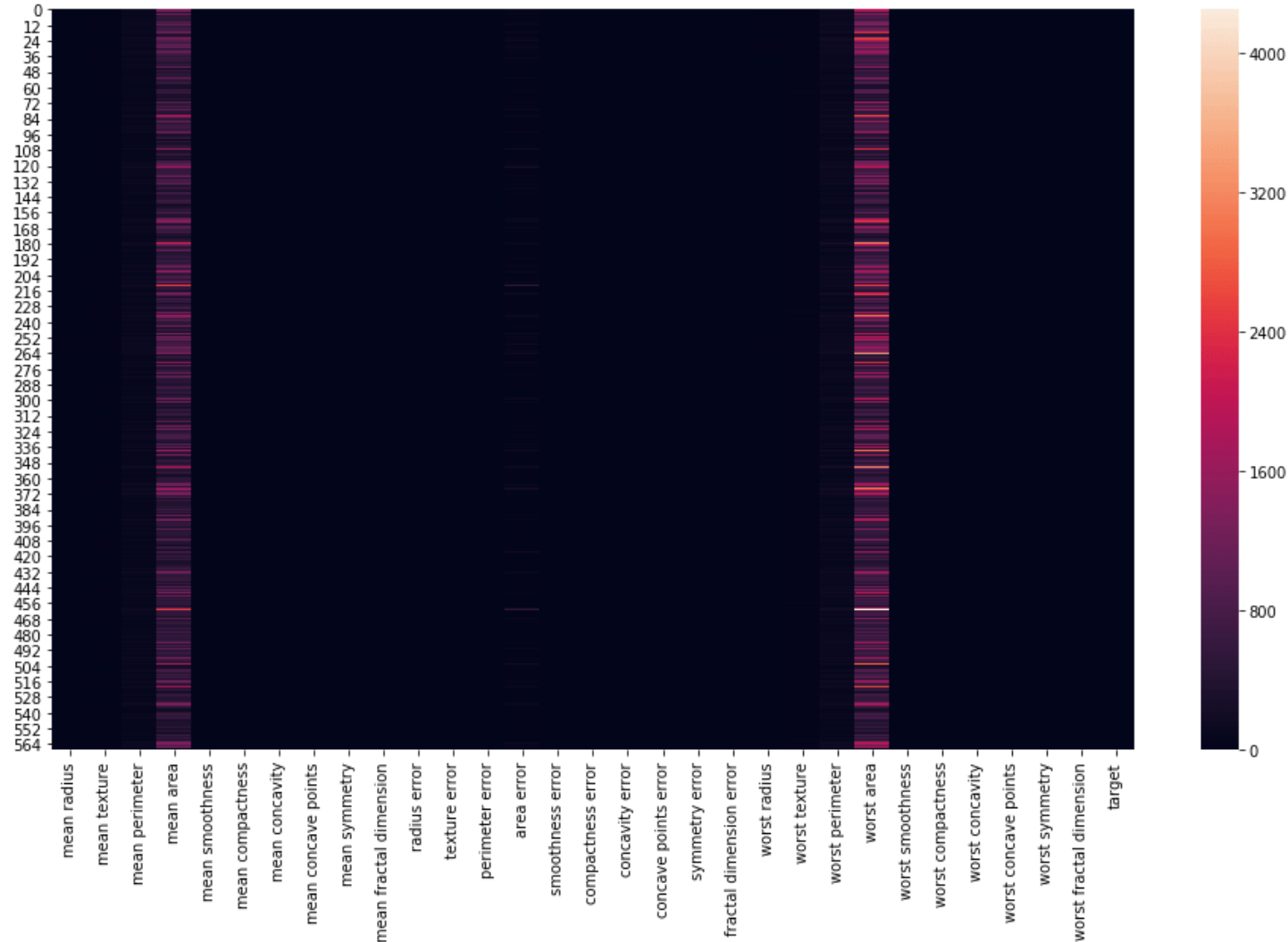
Heatmap

Heatmap of breast cancer DataFrame

In the below heatmap we can see the variety of different feature's value. The value of feature '**mean area**' and '**worst area**' are greater than other and 'mean perimeter', 'area error', and 'worst perimeter' value slightly less but greater than remaining features.

```
1 # heatmap of DataFrame
2 plt.figure(figsize=(16,9))
3 sns.heatmap(cancer_df)
```

Output >>>

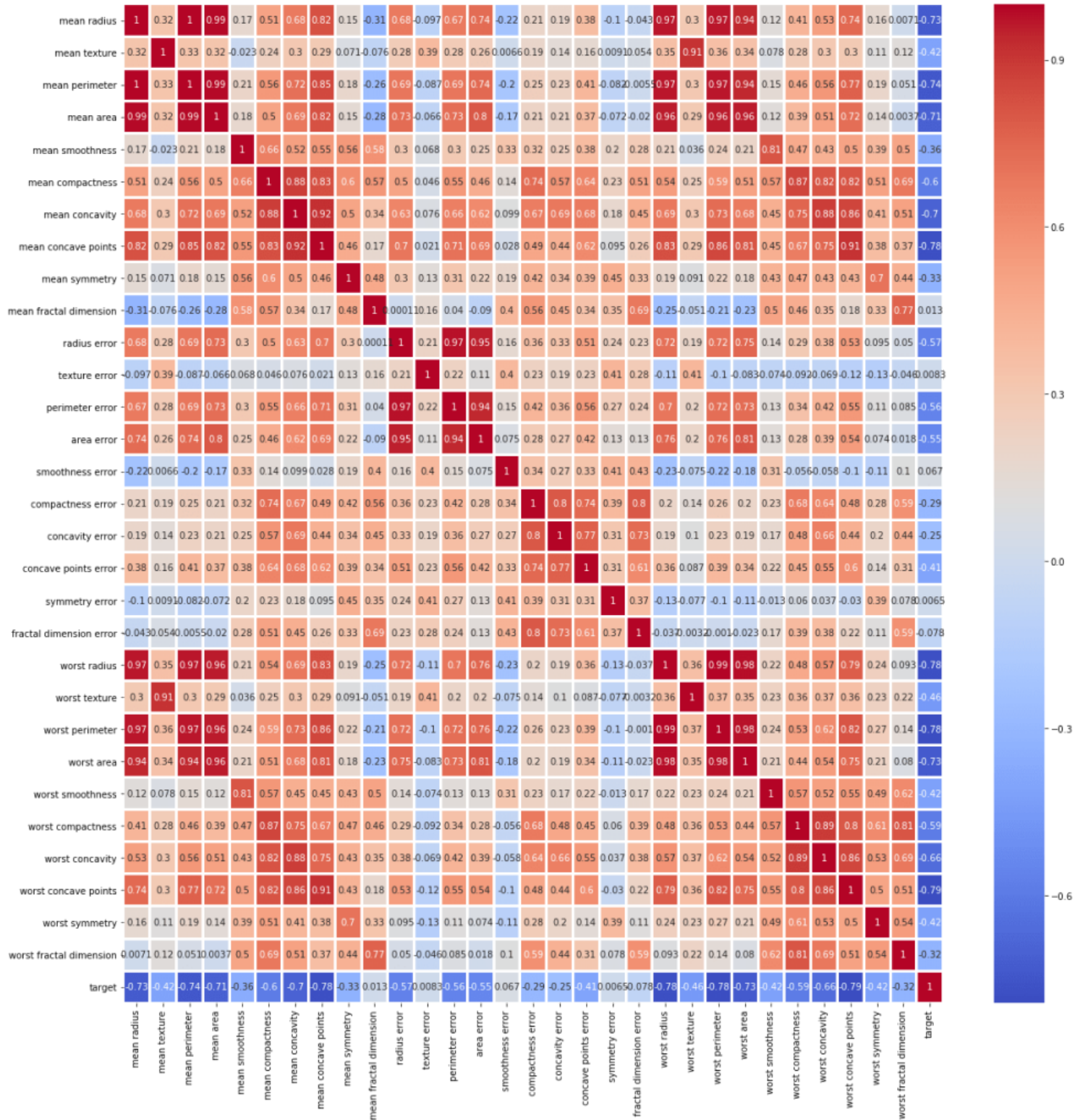


Heatmap of a correlation matrix

To find a correlation between each feature and target we visualize heatmap using the correlation matrix.

```
1 # Heatmap of Correlation matrix of breast cancer DataFrame
2 plt.figure(figsize=(20,20))
3 sns.heatmap(cancer_df.corr(), annot = True, cmap = 'coolwarm', linewidths=2)
```

Output >>>



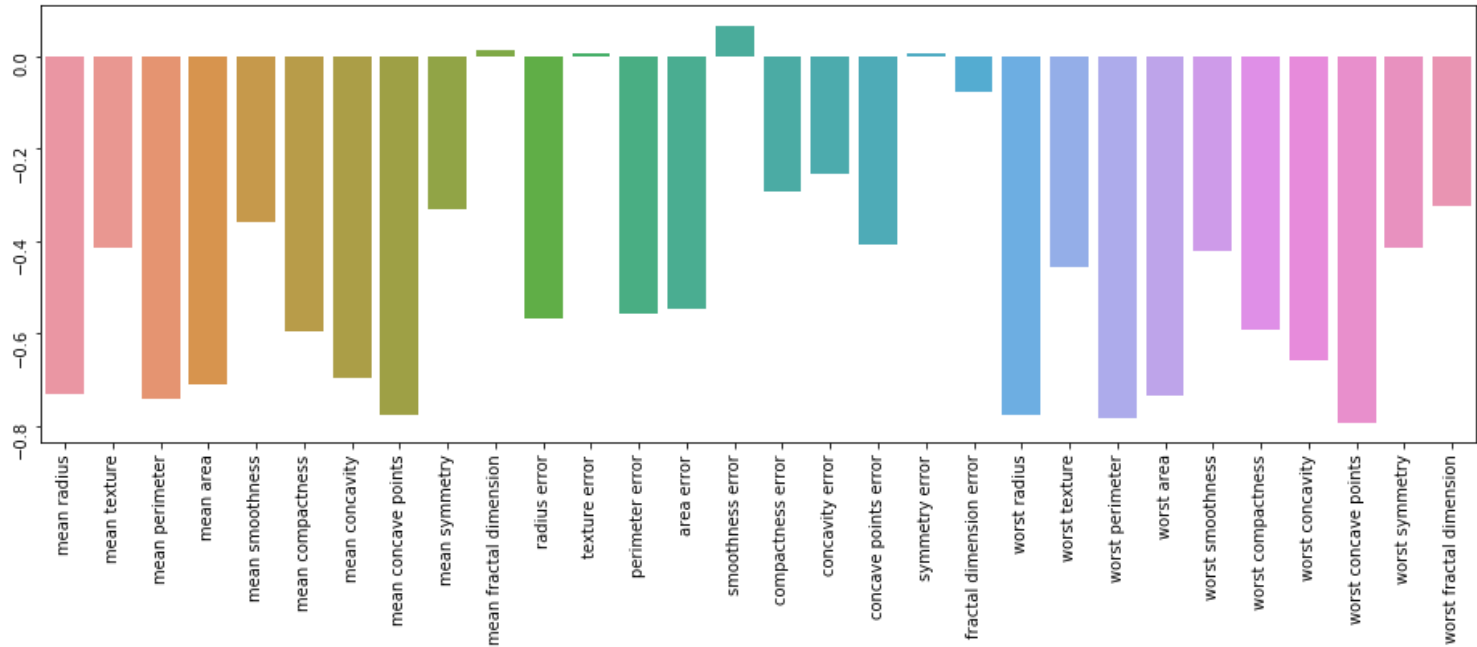
Correlation barplot

Taking the correlation of each feature with the target and the visualize barplot.

```
1 # create second DataFrame by dropping target
2 cancer_df2 = cancer_df.drop(['target'], axis = 1)
3 print("The shape of 'cancer_df2' is :", cancer_df2.shape)

1 # visualize correlation barplot
2 plt.figure(figsize = (16,5))
3 ax = sns.barplot(cancer_df2.corrwith(cancer_df.target).index, cancer_df2.corrwith(cancer_df.target))
4 ax.tick_params(labelrotation = 90)
```

Output >>>



In the above correlation barplot only feature 'smoothness error' is strongly positively correlated with the target than others. The features 'mean factor dimension', 'texture error', and 'symmetry error' are very less positive correlated and others remaining are strongly negatively correlated.

Data Preprocessing

Split DataFrame in train and test

```
1 | # input variable
2 | X = cancer_df.drop(['target'], axis = 1)
3 | X.head(6)
```

Output >>>

	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	worst area	worst smoothness	worst compactness	worst concavity	worst concave points	worst symmetry	worst fractal dimension
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	2019.0	0.1622	0.6656	0.7119	0.2654	0.4601	0.11890
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	1956.0	0.1238	0.1866	0.2416	0.1860	0.2750	0.08902
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	1709.0	0.1444	0.4245	0.4504	0.2430	0.3613	0.08758
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	567.7	0.2098	0.8663	0.6869	0.2575	0.6638	0.17300
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	1575.0	0.1374	0.2050	0.4000	0.1625	0.2364	0.07678
5	12.45	15.70	82.57	477.1	0.12780	0.17000	0.1578	0.08089	0.2087	0.07613	...	15.47	23.75	103.40	741.6	0.1791	0.5249	0.5355	0.1741	0.3985	0.12440

6 rows x 30 columns

```
1 | # output variable
2 | y = cancer_df['target']
3 | y.head(6)
```

Output >>>

```
1 | 0    0.0
2 | 1    0.0
3 | 2    0.0
4 | 3    0.0
5 | 4    0.0
6 | 5    0.0
7 | Name: target, dtype: float64
```

```
1 | # split dataset into train and test
2 | from sklearn.model_selection import train_test_split
3 | X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state= 5)
```

Feature Scaling

Converting different units and magnitude data in one unit.

```
1 | # Feature scaling
2 | from sklearn.preprocessing import StandardScaler
3 | sc = StandardScaler()
4 | X_train_sc = sc.fit_transform(X_train)
5 | X_test_sc = sc.transform(X_test)
```

Breast Cancer Detection Machine Learning Model Building

To clean data to build the ML model. But which Machine learning algorithm is best for the data we have to find. The output is a categorical format so will use supervised classification machine learning algorithms.

To build the best model, to train and test the dataset with multiple Machine Learning algorithms then we can find the best ML model. So let's try.

First, need to import the required packages.

```
1 | from sklearn.metrics import confusion_matrix, classification_report, accuracy_score
```

Support Vector Classifier

```
1 | # Support vector classifier
2 | from sklearn.svm import SVC
3 | svc_classifier = SVC()
4 | svc_classifier.fit(X_train, y_train)
5 | y_pred_scv = svc_classifier.predict(X_test)
6 | accuracy_score(y_test, y_pred_scv)
```

Output >>> 0.5789473684210527

```
1 | # Train with Standard scaled Data
2 | svc_classifier2 = SVC()
3 | svc_classifier2.fit(X_train_sc, y_train)
4 | y_pred_svc_sc = svc_classifier2.predict(X_test_sc)
5 | accuracy_score(y_test, y_pred_svc_sc)
```

Output >>> 0.9649122807017544

Logistic Regression

```
1 | # Logistic Regression
2 | from sklearn.linear_model import LogisticRegression
3 | lr_classifier = LogisticRegression(random_state = 51, penalty = 'l1')
4 | lr_classifier.fit(X_train, y_train)
5 | y_pred_lr = lr_classifier.predict(X_test)
6 | accuracy_score(y_test, y_pred_lr)
```

Output >>> 0.9736842105263158

```
1 | # Train with Standard scaled Data
2 | lr_classifier2 = LogisticRegression(random_state = 51, penalty = 'l1')
3 | lr_classifier2.fit(X_train_sc, y_train)
4 | y_pred_lr_sc = lr_classifier.predict(X_test_sc)
5 | accuracy_score(y_test, y_pred_lr_sc)
```

Output >>> 0.5526315789473685

K – Nearest Neighbor Classifier

```
1 | # K - Nearest Neighbor Classifier
2 | from sklearn.neighbors import KNeighborsClassifier
3 | knn_classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
4 | knn_classifier.fit(X_train, y_train)
5 | y_pred_knn = knn_classifier.predict(X_test)
6 | accuracy_score(y_test, y_pred_knn)
```

Output >>> 0.9385964912280702

```
1 | # Train with Standard scaled Data
2 | knn_classifier2 = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
3 | knn_classifier2.fit(X_train_sc, y_train)
4 | y_pred_knn_sc = knn_classifier.predict(X_test_sc)
5 | accuracy_score(y_test, y_pred_knn_sc)
```

Output >>> 0.5789473684210527

Naive Bayes Classifier

```
1 | # Naive Bayes Classifier
2 | from sklearn.naive_bayes import GaussianNB
3 | nb_classifier = GaussianNB()
4 | nb_classifier.fit(X_train, y_train)
5 | y_pred_nb = nb_classifier.predict(X_test)
6 | accuracy_score(y_test, y_pred_nb)
```

Output >>> 0.9473684210526315

```
1 | # Train with Standard scaled Data
2 | nb_classifier2 = GaussianNB()
3 | nb_classifier2.fit(X_train_sc, y_train)
4 | y_pred_nb_sc = nb_classifier2.predict(X_test_sc)
5 | accuracy_score(y_test, y_pred_nb_sc)
```

Output >>> 0.9385964912280702

Decision Tree Classifier

```
1 | # Decision Tree Classifier
2 | from sklearn.tree import DecisionTreeClassifier
3 | dt_classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 51)
4 | dt_classifier.fit(X_train, y_train)
5 | y_pred_dt = dt_classifier.predict(X_test)
6 | accuracy_score(y_test, y_pred_dt)
```

Output >>> 0.9473684210526315

```
1 # Train with Standard scaled Data
2 dt_classifier2 = DecisionTreeClassifier(criterion = 'entropy', random_state = 51)
3 dt_classifier2.fit(X_train_sc, y_train)
4 y_pred_dt_sc = dt_classifier.predict(X_test_sc)
5 accuracy_score(y_test, y_pred_dt_sc)
```

Output >>> 0.7543859649122807

Random Forest Classifier

```
1 # Random Forest Classifier
2 from sklearn.ensemble import RandomForestClassifier
3 rf_classifier = RandomForestClassifier(n_estimators = 20, criterion = 'entropy', random_state = 51)
4 rf_classifier.fit(X_train, y_train)
5 y_pred_rf = rf_classifier.predict(X_test)
6 accuracy_score(y_test, y_pred_rf)
```

Output >>> 0.9736842105263158

```
1 # Train with Standard scaled Data
2 rf_classifier2 = RandomForestClassifier(n_estimators = 20, criterion = 'entropy', random_state = 51)
3 rf_classifier2.fit(X_train_sc, y_train)
4 y_pred_rf_sc = rf_classifier.predict(X_test_sc)
5 accuracy_score(y_test, y_pred_rf_sc)
```

Output >>> 0.7543859649122807

Adaboost Classifier

```
1 # Adaboost Classifier
2 from sklearn.ensemble import AdaBoostClassifier
3 adb_classifier = AdaBoostClassifier(DecisionTreeClassifier(criterion = 'entropy', random_state = 200),
4                                     n_estimators=2000,
5                                     learning_rate=0.1,
6                                     algorithm='SAMME.R',
7                                     random_state=1,)
8 adb_classifier.fit(X_train, y_train)
9 y_pred_adb = adb_classifier.predict(X_test)
10 accuracy_score(y_test, y_pred_adb)
```

Output >>> 0.9473684210526315

```
1 # Train with Standard scaled Data
2 adb_classifier2 = AdaBoostClassifier(DecisionTreeClassifier(criterion = 'entropy', random_state = 200),
3                                     n_estimators=2000,
4                                     learning_rate=0.1,
5                                     algorithm='SAMME.R',
6                                     random_state=1,)
7 adb_classifier2.fit(X_train_sc, y_train)
8 y_pred_adb_sc = adb_classifier2.predict(X_test_sc)
9 accuracy_score(y_test, y_pred_adb_sc)
```

Output >>> 0.9473684210526315

XGBoost Classifier

```
1 # XGBoost Classifier
2 from xgboost import XGBClassifier
3 xgb_classifier = XGBClassifier()
4 xgb_classifier.fit(X_train, y_train)
5 y_pred_xgb = xgb_classifier.predict(X_test)
6 accuracy_score(y_test, y_pred_xgb)
```

Output >>> 0.9824561403508771

```
1 # Train with Standard scaled Data
2 xgb_classifier2 = XGBClassifier()
3 xgb_classifier2.fit(X_train_sc, y_train)
4 y_pred_xgb_sc = xgb_classifier2.predict(X_test_sc)
5 accuracy_score(y_test, y_pred_xgb_sc)
```

Output >>> 0.9824561403508771

XGBoost Parameter Tuning Randomized Search

```
1 # XGBoost classifier most required parameters
2 params={
3     "learning_rate" : [0.05, 0.10, 0.15, 0.20, 0.25, 0.30 ] ,
4     "max_depth" : [ 3, 4, 5, 6, 8, 10, 12, 15],
5     "min_child_weight" : [ 1, 3, 5, 7 ],
```

```
6 | "gamma"          : [ 0.0, 0.1, 0.2 , 0.3, 0.4 ],
7 | "colsample_bytree" : [ 0.3, 0.4, 0.5 , 0.7 ]
8 | }
```

```
1 | # Randomized Search
2 | from sklearn.model_selection import RandomizedSearchCV
3 | random_search = RandomizedSearchCV(xgb_classifier, param_distributions=params, scoring= 'roc_auc', n_jobs= -1, verbose=
4 | random_search.fit(X_train, y_train)
```

Output >>>

```
1 | RandomizedSearchCV(cv='warn', error_score='raise-deprecating',
2 |                   estimator=XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
3 |                   colsample_bynode=1, colsample_bytree=1, gamma=0, learning_rate=0.1,
4 |                   max_delta_step=0, max_depth=3, min_child_weight=1, missing=None,
5 |                   n_estimators=100, n_jobs=1, nthread=None,
6 |                   objective='binary:logistic', random_state=0, reg_alpha=0,
7 |                   reg_lambda=1, scale_pos_weight=1, seed=None, silent=None,
8 |                   subsample=1, verbosity=1),
9 |                   fit_params=None, iid='warn', n_iter=10, n_jobs=-1,
10 |                  param_distributions={'learning_rate': [0.05, 0.1, 0.15, 0.2, 0.25, 0.3], 'max_depth': [3, 4, 5, 6, 8, 10, 12
11 |                  pre_dispatch='2*n_jobs', random_state=None, refit=True,
12 |                  return_train_score='warn', scoring='roc_auc', verbose=3)
```

```
1 | random_search.best_params_
```

Output >>>

```
1 | {'min_child_weight': 1,
2 |  'max_depth': 3,
3 |  'learning_rate': 0.3,
4 |  'gamma': 0.4,
5 |  'colsample_bytree': 0.3}
```

```
1 | random_search.best_estimator_
```

Output >>>

```
1 | XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
2 |               colsample_bynode=1, colsample_bytree=0.3, gamma=0.4,
3 |               learning_rate=0.3, max_delta_step=0, max_depth=3,
4 |               min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
5 |               nthread=None, objective='binary:logistic', random_state=0,
6 |               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
7 |               silent=None, subsample=1, verbosity=1)
```

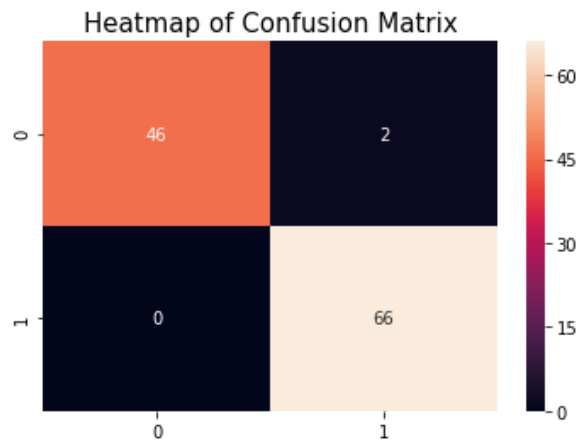
```
1 | # training XGBoost classifier with best parameters
2 | xgb_classifier_pt = XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
3 |               colsample_bynode=1, colsample_bytree=0.4, gamma=0.2,
4 |               learning_rate=0.1, max_delta_step=0, max_depth=15,
5 |               min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
6 |               nthread=None, objective='binary:logistic', random_state=0,
7 |               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
8 |               silent=None, subsample=1, verbosity=1)
9 |
10 | xgb_classifier_pt.fit(X_train, y_train)
11 | y_pred_xgb_pt = xgb_classifier_pt.predict(X_test)
```

```
1 | accuracy_score(y_test, y_pred_xgb_pt)
```

Output >>> 0.9824561403508771

Confusion Matrix

```
1 | cm = confusion_matrix(y_test, y_pred_xgb_pt)
2 | plt.title('Heatmap of Confusion Matrix', fontsize = 15)
3 | sns.heatmap(cm, annot = True)
4 | plt.show()
```



The model is giving 0% type II error and it is best.

Classification Report of Model

```
1 | print(classification_report(y_test, y_pred_xgb_pt))
```

Output >>>

	precision	recall	f1-score	support
0	0.0	1.00	0.96	48
1	1.0	0.97	1.00	66
micro avg	0.98	0.98	0.98	114
macro avg	0.99	0.98	0.98	114
weighted avg	0.98	0.98	0.98	114

Cross-validation of the ML model

To find the ML model is overfitted, under fitted or generalize doing cross-validation.

```
1 | # Cross validation
2 | from sklearn.model_selection import cross_val_score
3 | cross_validation = cross_val_score(estimator = xgb_model_pt2, X = X_train_sc, y = y_train, cv = 10)
4 | print("Cross validation of XGBoost model = ",cross_validation)
5 | print("Cross validation of XGBoost model (in mean) = ",cross_validation.mean())
6 | from sklearn.model_selection import cross_val_score
7 | cross_validation = cross_val_score(estimator = xgb_classifier_pt, X = X_train_sc,y = y_train, cv = 10)
8 | print("Cross validation accuracy of XGBoost model = ", cross_validation)
9 | print("\nCross validation mean accuracy of XGBoost model = ", cross_validation.mean())
```

Output >>>

1	Cross validation accuracy of XGBoost model =	[0.9787234 0.97826087 0.97826087 0.97826087 0.93333333 0.91111111
2	1.	1. 0.97777778 0.88888889]
3		
4	Cross validation mean accuracy of XGBoost model =	0.9624617124062083

The mean accuracy value of **cross-validation is 96.24%** and **XGBoost model accuracy is 98.24%**. It showing XGBoost is slightly overfitted but when training data will more it will generalized model.

Save the Machine Learning model

After completion of the Machine Learning project or building the ML model need to deploy in an application. To deploy the ML model need to save it first. To save the Machine Learning project we can use the **pickle** or **joblib** package.

```
1 | ## Pickle
2 | import pickle
3 |
4 | # save model
5 | pickle.dump(xgb_classifier_pt, open('breast_cancer_detector.pickle', 'wb'))
6 |
7 | # load model
8 | breast_cancer_detector_model = pickle.load(open('breast_cancer_detector.pickle', 'rb'))
9 |
10 | # predict the output
11 | y_pred = breast_cancer_detector_model.predict(X_test)
12 |
13 | # confusion matrix
14 | print('Confusion matrix of XGBoost model: \n',confusion_matrix(y_test, y_pred),'\n')
15 |
16 | # show the accuracy
17 | print('Accuracy of XGBoost model = ',accuracy_score(y_test, y_pred))
```

Output >>>

1	Confusion matrix of XGBoost model:
2	[[46 2]


```
3 | [ 0 66]]
4 |
5 | Accuracy of XGBoost model = 0.9824561403508771
```

I have completed the **Machine learning Project** successfully with **98.24%** accuracy which is great for '**Breast Cancer Detection using Machine learning**' project.