

predictor

April 29, 2024

```
[1]: import pandas as pd
import seaborn as sns
```

```
data = pd.read_csv("data.csv")
data.head()
```

```
[1]:      id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean \
0    842302        M      17.99      10.38      122.80      1001.0
1    842517        M      20.57      17.77      132.90      1326.0
2  84300903        M      19.69      21.25      130.00      1203.0
3  84348301        M      11.42      20.38       77.58       386.1
4  84358402        M      20.29      14.34      135.10      1297.0
```

```
      smoothness_mean  compactness_mean  concavity_mean  concave points_mean \
0          0.11840      0.27760      0.3001          0.14710
1          0.08474      0.07864      0.0869          0.07017
2          0.10960      0.15990      0.1974          0.12790
3          0.14250      0.28390      0.2414          0.10520
4          0.10030      0.13280      0.1980          0.10430
```

```
      ... texture_worst  perimeter_worst  area_worst  smoothness_worst \
0      ...      17.33      184.60      2019.0      0.1622
1      ...      23.41      158.80      1956.0      0.1238
2      ...      25.53      152.50      1709.0      0.1444
3      ...      26.50       98.87       567.7      0.2098
4      ...      16.67      152.20      1575.0      0.1374
```

```
      compactness_worst  concavity_worst  concave points_worst  symmetry_worst \
0          0.6656      0.7119          0.2654          0.4601
1          0.1866      0.2416          0.1860          0.2750
2          0.4245      0.4504          0.2430          0.3613
3          0.8663      0.6869          0.2575          0.6638
4          0.2050      0.4000          0.1625          0.2364
```

```
      fractal_dimension_worst  Unnamed: 32
0          0.11890      NaN
1          0.08902      NaN
2          0.08758      NaN
```

```

3          0.17300          NaN
4          0.07678          NaN

```

[5 rows x 33 columns]

```
[2]: data.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 569 entries, 0 to 568
Data columns (total 33 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   id                                     569 non-null    int64
1   diagnosis                             569 non-null    object
2   radius_mean                           569 non-null    float64
3   texture_mean                           569 non-null    float64
4   perimeter_mean                         569 non-null    float64
5   area_mean                             569 non-null    float64
6   smoothness_mean                       569 non-null    float64
7   compactness_mean                      569 non-null    float64
8   concavity_mean                        569 non-null    float64
9   concave points_mean                   569 non-null    float64
10  symmetry_mean                          569 non-null    float64
11  fractal_dimension_mean                 569 non-null    float64
12  radius_se                              569 non-null    float64
13  texture_se                             569 non-null    float64
14  perimeter_se                           569 non-null    float64
15  area_se                                569 non-null    float64
16  smoothness_se                          569 non-null    float64
17  compactness_se                         569 non-null    float64
18  concavity_se                           569 non-null    float64
19  concave points_se                      569 non-null    float64
20  symmetry_se                            569 non-null    float64
21  fractal_dimension_se                   569 non-null    float64
22  radius_worst                           569 non-null    float64
23  texture_worst                           569 non-null    float64
24  perimeter_worst                        569 non-null    float64
25  area_worst                             569 non-null    float64
26  smoothness_worst                       569 non-null    float64
27  compactness_worst                      569 non-null    float64
28  concavity_worst                        569 non-null    float64
29  concave points_worst                   569 non-null    float64
30  symmetry_worst                         569 non-null    float64
31  fractal_dimension_worst                 569 non-null    float64
32  Unnamed: 32                             0 non-null      float64
dtypes: float64(31), int64(1), object(1)
memory usage: 146.8+ KB

```

```
[3]: data.describe()
```

```
[3]:
```

	id	radius_mean	texture_mean	perimeter_mean	area_mean	\
count	5.690000e+02	569.000000	569.000000	569.000000	569.000000	
mean	3.037183e+07	14.127292	19.289649	91.969033	654.889104	
std	1.250206e+08	3.524049	4.301036	24.298981	351.914129	
min	8.670000e+03	6.981000	9.710000	43.790000	143.500000	
25%	8.692180e+05	11.700000	16.170000	75.170000	420.300000	
50%	9.060240e+05	13.370000	18.840000	86.240000	551.100000	
75%	8.813129e+06	15.780000	21.800000	104.100000	782.700000	
max	9.113205e+08	28.110000	39.280000	188.500000	2501.000000	

	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	\
count	569.000000	569.000000	569.000000	569.000000	
mean	0.096360	0.104341	0.088799	0.048919	
std	0.014064	0.052813	0.079720	0.038803	
min	0.052630	0.019380	0.000000	0.000000	
25%	0.086370	0.064920	0.029560	0.020310	
50%	0.095870	0.092630	0.061540	0.033500	
75%	0.105300	0.130400	0.130700	0.074000	
max	0.163400	0.345400	0.426800	0.201200	

	symmetry_mean	...	texture_worst	perimeter_worst	area_worst	\
count	569.000000	...	569.000000	569.000000	569.000000	
mean	0.181162	...	25.677223	107.261213	880.583128	
std	0.027414	...	6.146258	33.602542	569.356993	
min	0.106000	...	12.020000	50.410000	185.200000	
25%	0.161900	...	21.080000	84.110000	515.300000	
50%	0.179200	...	25.410000	97.660000	686.500000	
75%	0.195700	...	29.720000	125.400000	1084.000000	
max	0.304000	...	49.540000	251.200000	4254.000000	

	smoothness_worst	compactness_worst	concavity_worst	\
count	569.000000	569.000000	569.000000	
mean	0.132369	0.254265	0.272188	
std	0.022832	0.157336	0.208624	
min	0.071170	0.027290	0.000000	
25%	0.116600	0.147200	0.114500	
50%	0.131300	0.211900	0.226700	
75%	0.146000	0.339100	0.382900	
max	0.222600	1.058000	1.252000	

	concave points_worst	symmetry_worst	fractal_dimension_worst	\
count	569.000000	569.000000	569.000000	
mean	0.114606	0.290076	0.083946	
std	0.065732	0.061867	0.018061	
min	0.000000	0.156500	0.055040	

25%	0.064930	0.250400	0.071460
50%	0.099930	0.282200	0.080040
75%	0.161400	0.317900	0.092080
max	0.291000	0.663800	0.207500

```

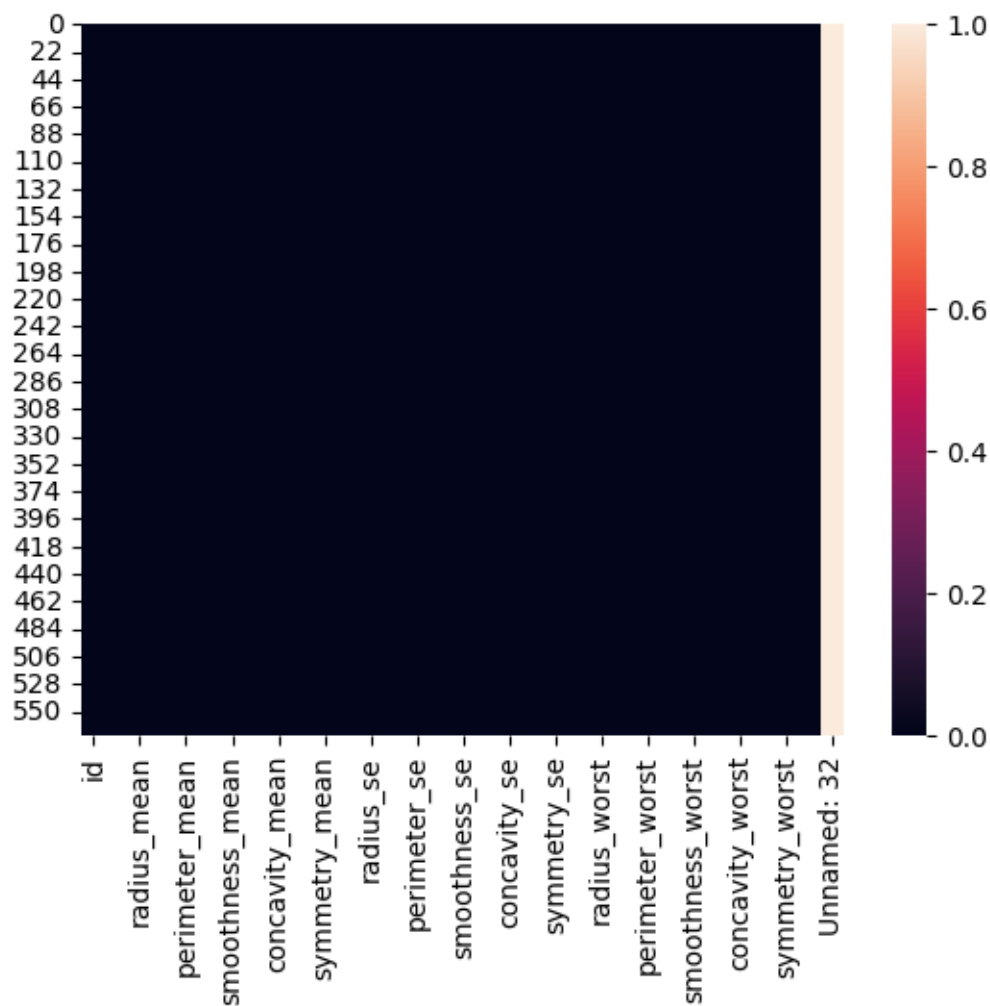
      Unnamed: 32
count      0.0
mean      NaN
std       NaN
min       NaN
25%       NaN
50%       NaN
75%       NaN
max       NaN

```

[8 rows x 32 columns]

```
[4]: sns.heatmap(data.isnull())
```

```
[4]: <Axes: >
```



```
[5]: data.drop(["Unnamed: 32", "id"], axis=1, inplace=True)
data.head()
```

```
[5]:
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	\
0	M	17.99	10.38	122.80	1001.0	
1	M	20.57	17.77	132.90	1326.0	
2	M	19.69	21.25	130.00	1203.0	
3	M	11.42	20.38	77.58	386.1	
4	M	20.29	14.34	135.10	1297.0	

	smoothness_mean	compactness_mean	concavity_mean	concave	points_mean	\
0	0.11840	0.27760	0.3001		0.14710	
1	0.08474	0.07864	0.0869		0.07017	
2	0.10960	0.15990	0.1974		0.12790	
3	0.14250	0.28390	0.2414		0.10520	

4	0.10030	0.13280	0.1980	0.10430
---	---------	---------	--------	---------

	symmetry_mean	...	radius_worst	texture_worst	perimeter_worst	\
0	0.2419	...	25.38	17.33	184.60	
1	0.1812	...	24.99	23.41	158.80	
2	0.2069	...	23.57	25.53	152.50	
3	0.2597	...	14.91	26.50	98.87	
4	0.1809	...	22.54	16.67	152.20	

	area_worst	smoothness_worst	compactness_worst	concavity_worst	\
0	2019.0	0.1622	0.6656	0.7119	
1	1956.0	0.1238	0.1866	0.2416	
2	1709.0	0.1444	0.4245	0.4504	
3	567.7	0.2098	0.8663	0.6869	
4	1575.0	0.1374	0.2050	0.4000	

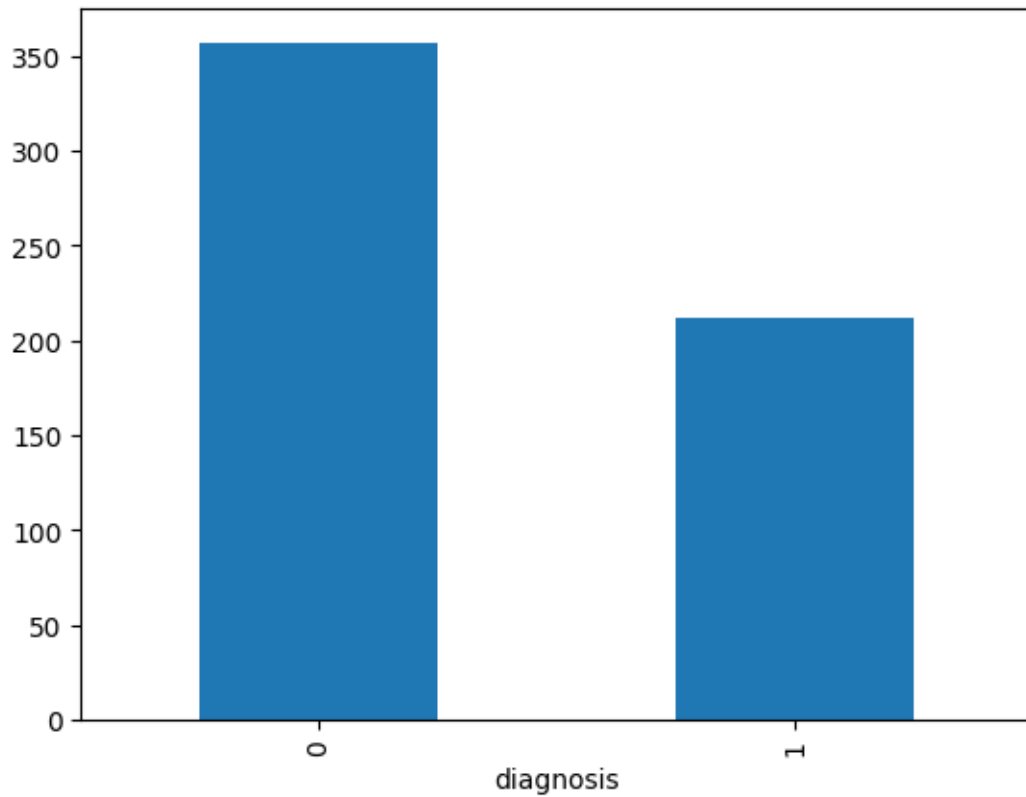
	concave	points_worst	symmetry_worst	fractal_dimension_worst
0		0.2654	0.4601	0.11890
1		0.1860	0.2750	0.08902
2		0.2430	0.3613	0.08758
3		0.2575	0.6638	0.17300
4		0.1625	0.2364	0.07678

[5 rows x 31 columns]

```
[6]: data.diagnosis = [1 if value == "M" else 0 for value in data.diagnosis]
```

```
[9]: data["diagnosis"].value_counts().plot(kind="bar")
```

```
[9]: <Axes: xlabel='diagnosis'>
```



```
[10]: # divide into target variable and predictors
y = data["diagnosis"] # target variable
x = data.drop(["diagnosis"], axis=1)
```

```
[11]: y
```

```
[11]: 0      1
      1      1
      2      1
      3      1
      4      1
      ..
     564      1
     565      1
     566      1
     567      1
     568      0
      Name: diagnosis, Length: 569, dtype: int64
```

```
[12]: x
```

```

[12]:      radius_mean  texture_mean  perimeter_mean  area_mean  smoothness_mean  \
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

      compactness_mean  concavity_mean  concave points_mean  symmetry_mean  \
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

      fractal_dimension_mean  ...  radius_worst  texture_worst  \
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

      perimeter_worst  area_worst  smoothness_worst  compactness_worst  \
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

	concavity_worst	concave points_worst	symmetry_worst	\
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	

	fractal_dimension_worst
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]

0.1 Normalize the data

```
[16]: from sklearn.preprocessing import StandardScaler

# Create a scaler object
scaler = StandardScaler()

# fit scaler to the data and transform the data
x_scaled = scaler.fit_transform(x)
```

```
[17]: x_scaled
```

```
[17]: array([[ 1.09706398, -2.07333501,  1.26993369, ...,  2.29607613,
          2.75062224,  1.93701461],
```

```
[ 1.82982061, -0.35363241, 1.68595471, ..., 1.0870843 ,
 -0.24388967, 0.28118999],
 [ 1.57988811, 0.45618695, 1.56650313, ..., 1.95500035,
 1.152255 , 0.20139121],
 ...,
 [ 0.70228425, 2.0455738 , 0.67267578, ..., 0.41406869,
 -1.10454895, -0.31840916],
 [ 1.83834103, 2.33645719, 1.98252415, ..., 2.28998549,
 1.91908301, 2.21963528],
 [-1.80840125, 1.22179204, -1.81438851, ..., -1.74506282,
 -0.04813821, -0.75120669]])
```

0.2 Split the data

```
[18]: from sklearn.model_selection import train_test_split

x_train, x_test, y_train, y_test = train_test_split(x_scaled, y, test_size=0.
↳30, random_state=42)
```

0.3 Train the model

```
[19]: from sklearn.linear_model import LogisticRegression

# Create the lr model
lr = LogisticRegression()

# Train the model on training data
lr.fit(x_train, y_train)

# Predict the target variable based on test data
y_predictions = lr.predict(x_test)
```

```
[20]: y_predictions
```

```
[20]: array([0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0,
 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 0,
 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 0,
 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1,
 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 0, 0, 0, 1, 0, 0,
 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 1,
 0, 0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0,
 0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0], dtype=int64)
```

1 Evaluation of the model

```
[22]: from sklearn.metrics import accuracy_score

accuracy = accuracy_score(y_test, y_predictions)
print(f"Accuracy: {accuracy: .2f}")
```

Accuracy: 0.98

```
[23]: from sklearn.metrics import classification_report
print(classification_report(y_test, y_predictions))
```

	precision	recall	f1-score	support
0	0.99	0.98	0.99	108
1	0.97	0.98	0.98	63
accuracy			0.98	171
macro avg	0.98	0.98	0.98	171
weighted avg	0.98	0.98	0.98	171

```
[ ]:
```