

# Predicting Breast Cancer in a patient

Problem statement: The objective of the dataset is to predict whether the cancer is benign or malignant

Data definition: The dataset consists of several predictor variables and one target variable. Diagnosis. The target variable has values 'Benign' and 'Malignant', where 'Benign' means that the cells are not harmful or there is no cancer and 'Malignant' means that the patient has cancer and the cells have a harmful effect

Table of contents: 1.Import packages 2.Read Data 3.Understand and prepare the data 3.1 Datatypes and dimensions 3.2 Statistical Summary 3.3 Label encoding for target variable 3.4 Missing data treatment 3.5 Visualization 3.6 Correlation 3.7 Feature selection 4.Ensemble learning techniques 4.1 Bagging meta-estimator

## 1.Import packages

```
In [2]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
# Suppress warnings
from warnings import filterwarnings
filterwarnings('ignore')
```

## 2.Read Data

```
In [3]: #read the data
data=pd.read_csv("cancer.csv")

#print first five rows of data
data.head()
```

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	...	texture_worst	perimeter_worst	area_worst	
Out[3]:	0	842302	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	17.33	184.60	2019.0
	1	842517	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	23.41	158.80	1956.0
	2	8430903	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	25.53	152.50	1709.0
	3	84348301	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	26.50	98.87	567.7
	4	84358402	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	16.67	152.20	1575.0

5 rows × 33 columns

## 3.Understand and prepare the Data

The process of data preparation entails cleansing, structuring and integrating data to make it ready for analysis.

### 3.1 Datatypes and Dimensions

```
In [4]: 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
```

From the above dataset, we see that 1. There is an id that cannot be used for classification 2. Unnamed:32 feature includes NaN so we can drop both the columns from the dataset

```
In [5]: data=data.drop(['id','Unnamed: 32'],axis=1)
data.head()
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	...	radius_worst	texture_worst	perimeter_worst
Out[5]:	0	M	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	25.38	17.33	
	1	M	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	24.99	23.41	
	2	M	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	23.57	25.53	
	3	M	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	14.91	26.50	
	4	M	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	22.54	16.67	

5 rows × 31 columns

```
In [7]: #shape of dataset
data.shape

Out[7]:
(569, 31)
```

```
In [8]: #plot the target
diagnosis.countplot(data['diagnosis'],label="count")
B=plt.ylabel("diagnosis").value_counts()
print("Number of B: ",B)
print("Number of M: ",M)
```

Number of B: 357  
Number of M: 212

### 3.2 Statistical Summary

In dataset there are numerical variables. describe() gives the statistical summary of numerical variables.

```
In [9]: data.describe()
```

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst
Out[9]:	count	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	569.000000	...	569.000000	...	569.000000	5
	mean	14.127322	19.289649	91.969033	654.889104	0.096360	0.104341	0.088799	0.048919	0.181162	...	0.062798	...	16.269180
	std	3.524049	4.301036	24.299861	351.914129	0.014064	0.052813	0.079720	0.038803	0.027414	...	0.007060	...	4.833242
	min	6.981000	9.710000	43.790000	143.500000	0.052630	0.019380	0.000000	0.000000	0.106000	...	0.049960	...	7.930000
	25%	11.700000	16.170000	75.170000	420.300000	0.086370	0.064620	0.029560	0.020310	0.161900	...	0.057700	...	13.010000
	50%	13.370000	18.840000	86.240000	551.300000	0.095870	0.092630	0.061540	0.033500	0.179200	...	0.061540	...	14.970000
	75%	15.780000	21.800000	104.100000	782.700000	0.105300	0.134040	0.130700	0.074000	0.195700	...	0.066120	...	18.780000
	max	28.110000	39.280000	188.500000	2501.000000	0.163400	0.345400	0.426800	0.201200	0.304000	...	0.097440	...	36.040000

8 rows × 30 columns

### 3.3 Label encoding for Target variable

Converting target variable into numeric

```
In [10]: #replace M with zero and B with one
data['diagnosis']=data['diagnosis'].replace('M',0)
data['diagnosis']=data['diagnosis'].replace('B',1)
data.head()
```

	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	...	radius_worst	texture_worst	perimeter_worst
Out[10]:	0	0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.3001	0.14710	...	25.38	17.33	
	1	0	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.0869	0.07017	...	24.99	23.41	
	2	0	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.1974	0.12790	...	23.57	25.53	
	3	0	11.42	20.38	77.58	386.1	0.14250	0.28390	0.2414	0.10520	...	14.91	26.50	
	4	0	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.1980	0.10430	...	22.54	16.67	

5 rows × 31 columns

### 3.4 Missing Data treatment

First we have to find whether there is any null values in our given data set. If there is null values we have to handle it properly , otherwise it effects our prediction

```
In [11]: #get the count of missing values
data.isnull().sum()
```

Out[11]:	diagnosis	0
	radius_mean	0
	texture_mean	0
	perimeter_mean	0
	area_mean	0
	smoothness_mean	0
	compactness_mean	0
	concavity_mean	0
	concave points_mean	0
	symmetry_mean	0
	fractal_dimension_mean	0
	radius_se	0
	texture_se	0
	perimeter_se	0
	area_se	0
	smoothness_se	0
	compactness_se	0
	concavity_se	0
	concave points_se	0
	symmetry_se	0
	fractal_dimension_se	0
	radius_worst	0
	texture_worst	0
	perimeter_worst	0
	area_worst	0
	smoothness_worst	0
	compactness_worst	0
	concavity_worst	0
	concave points_worst	0
	symmetry_worst	0
	fractal_dimension_worst	0
	dtype:	int64

There are no null values present in the dataset

### 3.5 Visualization

```
In [12]: hist_fig=data.hist(figsize=(18,18))
```

### Correlation

```
In [13]: target=data['diagnosis']
corr=target.corr()
corr
```

	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean	concavity_mean	concave points_mean	symmetry_mean	fractal_dimension_mean	...	radius_worst	texture_worst	perimeter_worst	area_worst
Out[13]:	radius_mean	1.000000	0.323782	0.997855	0.987357	0.170581	0.506124	0.676764	0.822529	0.147741	...	-0.311631	...		
	texture_mean	0.323782	1.000000	0.239533	0.321086	-0.023389	0.236702	0.302418	0.293464	0.071401	...	-0.076437	...		
	perimeter_mean	0.997855	0.329533	1.000000	0.986507	0.207278	0.556936	0.716136	0.850977	0.183027	...	-0.261477	...		
	area_mean	0.987357	0.321086	0.986507	1.000000	0.177028	0.498502	0.685983	0.823269	0.151293	...	-0.283110	...		
	smoothness_mean	0.170581	-0.023389	0.207278	0.177028	1.000000	0.659123	0.521984	0.535995	0.557775	...	0.584792	...		
	compactness_mean	0.506124	0.236702	0.556936	0.498502	0.659123	1.000000	0.883121	0.831135	0.602641	...	0.565369	...		
	concavity_mean	0.676764	0.302418	0.716136	0.685983	0.521984	0.883121	1.000000	0.921391	0.500667	...	0.376783	...		
	concave points_mean	0.822529	0.293464	0.850977	0.823269	0.535995	0.831125	0.921391	1.000000	0.462497	...	0.166917	...		
	symmetry_mean	0.147741	0.071401	0.183027	0.151293	0.557775	0.602641	0.500667	0.462497	1.000000	...	0.479921	...		
	fractal_dimension_mean	-0.311631	-0.076437	-0.261477	-0.283110	0.584792	0.565369	0.336783	0.166917	0.479921	...	1.000000	...		
	radius_se	0.679090	0.275869	0.691765	0.732562	0.301467	0.497473	0.631925	0.680550	0.303379	...	0.001011	...		
	texture_se	-0.097317	0.386358	-0.086761	-0.062280	0.068406	0.046205	0.076218	0.021480	0.128055	...	0.164174	...		
	perimeter_se	0.674172	0.281673	0.693135	0.726628	0.260992	0.548905	0.660391	0.710650	0.313893	...	0.039830	...		
	area_se	0.735864	0.258945	0.744983	0.800086	0.246552	0.456953	0.617427	0.690299	0.223970	...	-0.000170	...		
	smoothness_se	-0.222600	0.006614	-0.202094	-0.166777	0.332375	0.135299	0.090564	0.027653	0.187321	...	0.401964	...		
	compactness_se	0.206000	0.191975	0.250744	0.212583	0.318943	0.738722	0.670279	0.490424	0.421659	...	0.559837	...		
	concavity_se	0.194204	0.143293	0.228082	0.207660	0.248396	0.570517	0.691270	0.439167	0.342627	...	0.446630	...		
	concave points_se	0.736169	0.133951	0.803261	0.727320	0.300676	0.442262	0.683290	0.615634	0.393298	...	0.341198	...		
	symmetry_se	-0.104321	0.009127	-0.081629	-0.072497	0.200774	0.229977	0.178009	0.095351	0.449137	...	0.345007	...		
	fractal_dimension_se	-0.042641	0.054458	-0.021629	-0.019887	0.293607	0.259318	0.448301	0.257584	0.449137	...	0.686132	...		
	radius_worst	0.969639	0.352973	0.969476	0.962746	0.213120	0.535315	0.688236	0.830318	0.185758	...	-0.053691	...		
	texture_worst	0.297098	0.912045	0.300306	0.287489	0.036072	0.248133	0.299679	0.292752	0.090651	...	-0.021969	...		
	perimeter_worst	0.965137	0.358040	0.970387	0.959320	0.238853	0.590210	0.729565	0.855923	0.210169	...	-0.205151	...		
	area_worst	0.941062	0.343546	0.941559	0.959213	0.206718	0.509604	0.675887	0.809630	0.177193	...	-0.231954	...		
	smoothness_worst	0.119616	0.073503	0.164590	0.125253	0.865324	0.565541	0.448822	0.452753	0.426675	...	0.504942	...		
	compactness_worst	0.434643	0.277830	0.455774	0.390410	0.472468	0.865809	0.754968	0.667454	0.437320	...	0.468798	...		
	concavity_worst	0.526911	0.301025	0.563879	0.512606	0.434626	0.816275	0.884103	0.782399	0.433721	...	0.346234	...		
	concave points_worst	0.744214	0.295316	0.771241	0.722017	0.303093	0.815573	0.861323	0.910155	0.430287	...	0.173255	...		
	symmetry_worst	0.163953	0.105608	0.189115	0.143570	0.394309	0.510223	0.400464	0.375744	0.698626	...	0.334019	...		
	fractal_dimension_worst	0.070066	0.119205	0.050109	0.003738	0.499316	0.687382	0.514930	0.366661	0.438413	...	0.767297	...		

30 rows × 30 columns

```
In [23]: #correlation map
sns.heatmap(target.corr(),cmap='coolwarm')
plt.show()
```