# Cancer detection

November 12, 2024

## 0.1 Loading of modules

```python
[28]: import pandas as pd
      import numpy as np
      import seaborn as sns
      import matplotlib.pyplot as plt

      from sklearn.preprocessing import MinMaxScaler,StandardScaler
```

## 0.2 Load data

```python
[3]: df = pd.read_csv('Cancer_Data.csv',index_col=None).iloc[:,:-1]
     df.head()
```

```
[3]:          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

        …  radius_worst  texture_worst  perimeter_worst  area_worst  \
     0  …         25.38          17.33           184.60      2019.0
     1  …         24.99          23.41           158.80      1956.0
     2  …         23.57          25.53           152.50      1709.0
     3  …         14.91          26.50            98.87       567.7
     4  …         22.54          16.67           152.20      1575.0

        smoothness_worst  compactness_worst  concavity_worst  concave points_worst  \
     0            0.1622             0.6656           0.7119                0.2654
     1            0.1238             0.1866           0.2416                0.1860
```

```
2            0.1444            0.4245            0.4504            0.2430
3            0.2098            0.8663            0.6869            0.2575
4            0.1374            0.2050            0.4000            0.1625

   symmetry_worst  fractal_dimension_worst
0          0.4601                  0.11890
1          0.2750                  0.08902
2          0.3613                  0.08758
3          0.6638                  0.17300
4          0.2364                  0.07678

[5 rows x 32 columns]
```

[26]: `df.columns`

[26]: 
```
Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
       'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

## 0.3 Dataset Description

[50]: `df.describe()`

[50]:
```
                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  …  radius_worst  texture_worst  perimeter_worst  \
count      569.000000  …    569.000000     569.000000       569.000000
mean         0.181162  …     16.269190      25.677223       107.261213
std          0.027414  …      4.833242       6.146258        33.602542
min          0.106000  …      7.930000      12.020000        50.410000
25%          0.161900  …     13.010000      21.080000        84.110000
50%          0.179200  …     14.970000      25.410000        97.660000
75%          0.195700  …     18.790000      29.720000       125.400000
max          0.304000  …     36.040000      49.540000       251.200000

         area_worst  smoothness_worst  compactness_worst  concavity_worst  \
count    569.000000        569.000000         569.000000       569.000000
mean     880.583128          0.132369           0.254265         0.272188
std      569.356993          0.022832           0.157336         0.208624
min      185.200000          0.071170           0.027290         0.000000
25%      515.300000          0.116600           0.147200         0.114500
50%      686.500000          0.131300           0.211900         0.226700
75%     1084.000000          0.146000           0.339100         0.382900
max     4254.000000          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

[8 rows x 31 columns]
```

## 0.4 Different category of diagnosis

```
[4]: df.diagnosis.unique()
```

```
[4]: array(['M', 'B'], dtype=object)
```

## 0.5 Checking of Data Imbalancement

```
[17]: df[df.diagnosis=='M'].count()
```

```
[17]: id                        212
      diagnosis                 212
      radius_mean               212
      texture_mean              212
      perimeter_mean            212
      area_mean                 212
      smoothness_mean           212
      compactness_mean          212
      concavity_mean            212
      concave points_mean       212
      symmetry_mean             212
      fractal_dimension_mean    212
      radius_se                 212
      texture_se                212
      perimeter_se              212
      area_se                   212
      smoothness_se             212
      compactness_se            212
      concavity_se              212
      concave points_se         212
      symmetry_se               212
      fractal_dimension_se      212
      radius_worst              212
      texture_worst             212
      perimeter_worst           212
      area_worst                212
      smoothness_worst          212
      compactness_worst         212
      concavity_worst           212
      concave points_worst      212
      symmetry_worst            212
      fractal_dimension_worst   212
      dtype: int64
```

```
[18]: df[df.diagnosis=='B'].count()
```

```
[18]: id                        357
      diagnosis                 357
      radius_mean               357
      texture_mean              357
      perimeter_mean            357
      area_mean                 357
      smoothness_mean           357
      compactness_mean          357
      concavity_mean            357
      concave points_mean       357
      symmetry_mean             357
```

```
fractal_dimension_mean      357
radius_se                   357
texture_se                  357
perimeter_se                357
area_se                     357
smoothness_se               357
compactness_se              357
concavity_se                357
concave points_se           357
symmetry_se                 357
fractal_dimension_se        357
radius_worst                357
texture_worst               357
perimeter_worst             357
area_worst                  357
smoothness_worst            357
compactness_worst           357
concavity_worst             357
concave points_worst        357
symmetry_worst              357
fractal_dimension_worst     357
dtype: int64
```

## 0.6   Checking for null values if any

```
[22]: df[df.isna()==True].value_counts()
```

```
[22]: Series([], Name: count, dtype: int64)
```

## 0.7   Normalizing the data

```
[36]: scale = StandardScaler()
      scaled_data = scale.fit_transform(df.iloc[:,2:])
```

```
[47]: df.describe()
```

```
[47]:                  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  …  radius_worst  texture_worst  perimeter_worst  \
count     569.000000  …    569.000000     569.000000       569.000000
mean        0.181162  …     16.269190      25.677223       107.261213
std         0.027414  …      4.833242       6.146258        33.602542
min         0.106000  …      7.930000      12.020000        50.410000
25%         0.161900  …     13.010000      21.080000        84.110000
50%         0.179200  …     14.970000      25.410000        97.660000
75%         0.195700  …     18.790000      29.720000       125.400000
max         0.304000  …     36.040000      49.540000       251.200000

         area_worst  smoothness_worst  compactness_worst  concavity_worst  \
count    569.000000        569.000000         569.000000       569.000000
mean     880.583128          0.132369           0.254265         0.272188
std      569.356993          0.022832           0.157336         0.208624
min      185.200000          0.071170           0.027290         0.000000
25%      515.300000          0.116600           0.147200         0.114500
50%      686.500000          0.131300           0.211900         0.226700
75%     1084.000000          0.146000           0.339100         0.382900
max     4254.000000          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

[8 rows x 31 columns]
```

## 0.8 Figuring out the outliers

```
[46]: plt.figure(figsize=(15,10))
      sns.boxplot(scaled_data)
      x = [i for i in range(-1,len(df.columns))]
      y = [3.5 for i in range(len(df.columns)+1)]
      plt.plot(x,y,marker = '*',label = 'threshold')
      plt.legend()
```

[46]: <matplotlib.legend.Legend at 0x7cea4cade9b0>