

```
[1]: import pandas as pd  
      import numpy as np  
      import sklearn as sklearn
```

1. Loading and Preprocessing (2 marks)etc

- Load the breast cancer dataset from sklearn.
 - Preprocess the data to handle any missing values and perform necessary feature scaling.
 - Explain the preprocessing steps you performed and justify why they are necessary for this dataset.

```
[5]: df=pd.DataFrame(data.data,columns=data.feature_names)
df["target"]=data.target
```

[5]:	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
0	17.99	10.38	122.80	1001.0	0.11840	0.27760	0.30010	0.14710	0.2419	0.07871	...	17.33	184.60	2019.0	0.16220	0.66560
1	20.57	17.77	132.90	1326.0	0.08474	0.07864	0.06890	0.07017	0.1812	0.05667	...	23.41	158.80	1956.0	0.12380	0.18660
2	19.69	21.25	130.00	1203.0	0.10960	0.15990	0.19740	0.12790	0.2069	0.05999	...	25.53	152.50	1709.0	0.14440	0.42450
3	11.42	20.38	77.58	386.1	0.14250	0.28390	0.24140	0.10520	0.2597	0.09744	...	26.50	98.87	567.7	0.20980	0.86360
4	20.29	14.34	135.10	1297.0	0.10030	0.13280	0.19800	0.10430	0.1809	0.05883	...	16.67	152.20	1575.0	0.13740	0.20500
...
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
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
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
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
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

569 rows × 31 columns

```
[6]: # To get the basic information about the dataset like no: of rows,no: of columns, datatypes etc  
dt.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  target          569 non-null    int32 
dtypes: float64(30), int32(1)
memory usage: 18.5 KB
```

```
[7]: # To find the no of duplicated values  
df.duplicated().sum()
```

```
[9]: # To check is there any null values or missing values in the dataset.
```

```
[9]: 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  
target                  0  
dtype: int64
```

This dataset set do not contain any missing values.

```
[11]: # To analyse the statistical measures of the dataset.
```

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

	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
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
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	880.58312
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	569.35699
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	185.20000
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	515.30000
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	686.50000
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	1084.00000
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	4254.00000

8 rows × 31 columns

In some of the columns, there is major difference between the mean and median values. This may be the indication of outliers in the dataset.

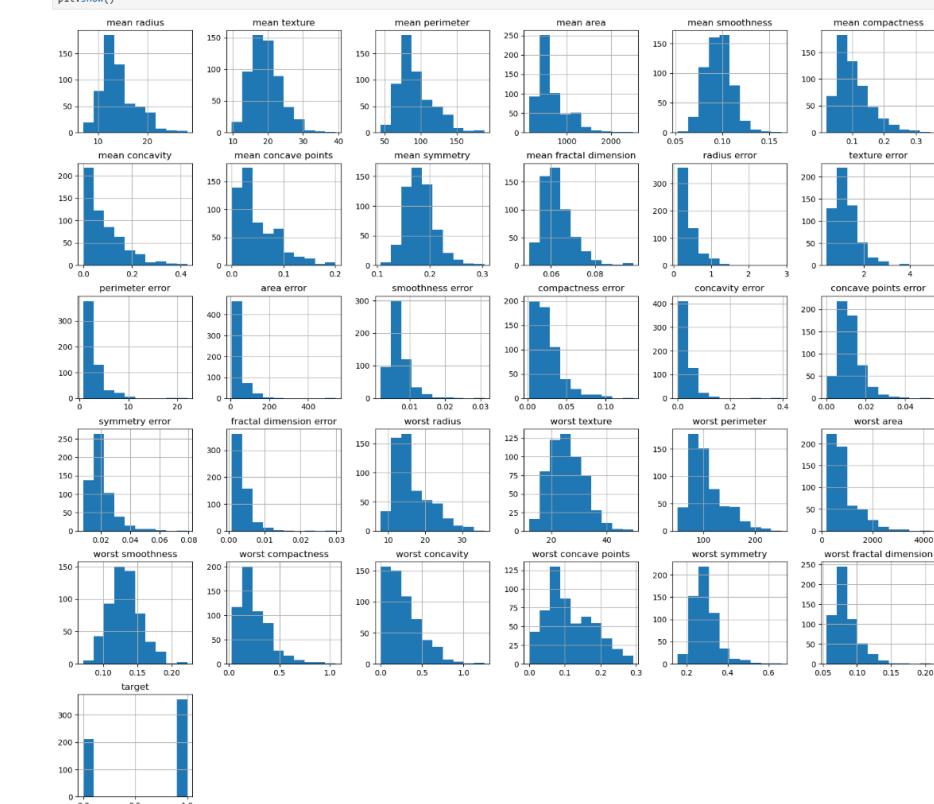
```
[13]: # To check the presence of outliers visually
```

```
[13]: import seaborn as sns  
import matplotlib.pyplot as plt  
numericals=df.select_dtypes("number")  
for column in numericals:  
    numericals.boxplot(column)  
plt.show()
```



we can see few outliers in the boxplots

```
[15]: df.hist(figsize=(20,18))  
plt.show()
```



From the histogram, we can see that the data does not have a uniform bell shape. This indicates the presence of outliers.

```
[17]: # To check the skewness of the data  
skewness=df.skew()  
skewness
```

mean radius	0.042388
mean texture	0.650459
mean perimeter	0.299659
mean area	1.645732
mean smoothness	0.456324
mean compactness	1.199123
mean concavity	1.401180
mean concave points	1.171180
mean symmetry	0.725609
mean fractal dimension	1.304489
radius error	3.088612
texture error	1.646444
perimeter error	3.443615
area error	5.447186
smoothness error	2.314450
compactness error	1.902221
concavity error	5.118463
concave points error	1.444678
symmetry error	2.195133
fractal dimension error	3.923969
worst radius	1.103115
worst texture	0.498321
worst perimeter	1.128164
...

```
worst area      1.859513
worst smoothness   0.415426
worst compactness   1.473555
worst concavity    1.150237
worst concave points  0.493616
worst symmetry     1.433928
worst fractal dimension  1.662579
target          -0.528461
dtype: float64
```

From the skewness result, we can see most of the columns are positively skewed.

```
[19]: # to find the outliers using IQR method.
def find_outliers_iqr(df):
    outliers= []
    numericals=df.select_dtypes(include=["number"])
    for column in numericals:
        Q1= numericals[column].quantile(0.25)
        Q3= numericals[column].quantile(0.75)
        IQR=Q3 - Q1
        lower_whisker=Q1 - 1.5*IQR
        upper_whisker=Q3 + 1.5*IQR
        outliers[column] = df[(numericals[column] < lower_whisker) | (numericals[column] > upper_whisker)]
    return outliers
outliers_dict=find_outliers_iqr(df)
for column,outliers in outliers_dict.items():
    print("outliers in",column,":",outliers)
```

	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
82	25.22	24.97	171.5	1079.0	0.106386	0.2665	0.3339	0.18450	0.1829	0.2665	25.27	24.97	152.0	1593.0	0.13269	0.144769
108	23.27	19.67	152.0	1054.0	0.13269	0.2665	0.3339	0.18450	0.1829	0.2665	24.25	20.28	166.2	1761.0	0.144769	0.144769
122	23.27	21.87	152.1	1696.0	0.08439	0.2665	0.3339	0.18450	0.1829	0.2665	23.27	22.04	152.1	2250.0	0.10940	0.10940
164	23.27	21.87	182.1	1685.0	0.11410	0.2665	0.3339	0.18450	0.1829	0.2665	23.27	21.87	182.1	2499.0	0.11420	0.11420
188	23.27	21.87	188.5	1685.0	0.11420	0.2665	0.3339	0.18450	0.1829	0.2665	23.27	21.87	188.5	1682.0	0.09342	0.09342
202	23.29	26.67	158.9	1685.0	0.11420	0.2665	0.3339	0.18450	0.1829	0.2665	23.29	26.67	158.9	1685.0	0.11420	0.11420
212	28.11	18.47	188.5	2499.0	0.11420	0.2665	0.3339	0.18450	0.1829	0.2665	28.11	18.47	188.5	1682.0	0.09342	0.09342
236	23.21	26.97	153.5	1670.0	0.09509	0.2665	0.3339	0.18450	0.1829	0.2665	23.21	26.97	153.5	1670.0	0.09509	0.09509
339	23.51	24.27	155.1	1747.0	0.10659	0.2665	0.3339	0.18450	0.1829	0.2665	23.51	24.27	155.1	1747.0	0.10659	0.10659
352	25.73	17.46	174.2	2010.0	0.11490	0.2665	0.3339	0.18450	0.1829	0.2665	22.01	21.90	147.2	1482.0	0.10653	0.10653
369	27.42	26.27	186.9	2501.0	0.10840	0.2665	0.3339	0.18450	0.1829	0.2665	27.42	26.27	186.9	2501.0	0.10840	0.10840
461	23.09	19.83	152.1	1682.0	0.09342	0.2665	0.3339	0.18450	0.1829	0.2665	23.09	19.83	152.1	1682.0	0.09342	0.09342
503	24.63	21.60	165.5	1841.0	0.10300	0.2665	0.3339	0.18450	0.1829	0.2665	24.63	21.60	165.5	1841.0	0.10300	0.10300

mean compactness mean concavity mean concave points mean symmetry \\\n82 0.2665 0.3339 0.18450 0.1829

```
[20]: #To cap the outliers
def cap_outliers(df):
    df_capped=df.copy()
    numericals=df.select_dtypes(include=["number"])

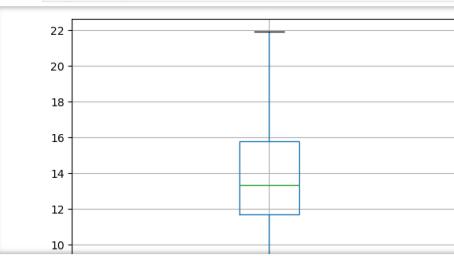
    for column in numericals.columns:
        Q1 = numericals[column].quantile(0.25)
        Q3 = numericals[column].quantile(0.75)
        IQR = Q3 - Q1
        lower_whisker = Q1 - 1.5 * IQR
        upper_whisker = Q3 + 1.5 * IQR
        df_capped[column]=df[column].apply(lambda x:lower_whisker if x < lower_whisker else upper_whisker if x > upper_whisker else x)

    return df_capped
df_capped=cap_outliers(df)
df_capped
```

	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
0	17.99	10.38	122.80	1001.0	0.118400	0.22862	0.28241	0.14710	0.2419	0.07871	17.33	184.60	1937.05	0.16220	0.62695	
1	20.57	17.77	132.90	1326.0	0.084740	0.07864	0.08690	0.07017	0.1812	0.05667	23.41	158.80	1937.05	0.12380	0.18660	
2	19.69	21.25	130.00	1203.0	0.109600	0.15990	0.19740	0.12790	0.2069	0.05999	25.53	152.50	1709.00	0.14440	0.42450	
3	11.42	20.38	77.58	386.1	0.133695	0.22862	0.24140	0.10520	0.2464	0.07875	26.50	98.87	567.70	0.19010	0.62695	
4	20.29	14.34	135.10	1297.0	0.100300	0.13280	0.19800	0.10430	0.1809	0.05883	16.67	152.20	1575.00	0.13740	0.20500	
...	
564	21.56	22.39	142.00	1326.3	0.111000	0.11590	0.24390	0.13890	0.1726	0.05623	26.40	166.10	1937.05	0.14100	0.21130	
565	20.13	28.25	131.20	1261.0	0.097800	0.10340	0.14400	0.09791	0.1752	0.05533	38.25	155.00	1731.00	0.11660	0.19220	
566	16.60	28.08	108.30	858.1	0.084550	0.10230	0.09251	0.05302	0.1590	0.05648	34.12	126.70	1124.00	0.11390	0.30940	
567	20.60	29.33	140.10	1265.0	0.117800	0.22662	0.28241	0.15200	0.2397	0.07016	39.42	184.60	1821.00	0.16500	0.62695	
568	7.76	24.54	47.92	181.0	0.059795	0.04362	0.00000	0.00000	0.1587	0.05884	30.37	59.16	268.60	0.08996	0.06444	

569 rows × 31 columns

```
[21]: # checking the presence of outliers after capping the outliers.
numericals=df_capped.select_dtypes("number")
for column in numericals:
    numericals.boxplot(column)
    plt.show()
```



outliers capped successfully

```
[23]: #checking the skewness of the data after capping the outliers
skewness=df_capped.skew()
```

```
[23]: mean radius      0.655953
mean texture      0.449700
mean perimeter    0.701811
mean area         0.952894
mean smoothness   0.257712
mean compactness   0.826755
mean concavity    1.023859
mean concave points  1.004049
mean symmetry     0.493621
mean fractal dimension  0.682430
radius error       0.259301
texture error      0.749987
perimeter error    1.034309
area error         1.130940
smoothness error   0.780923
compactness error  0.990285
concavity error    0.916740
concave points error  0.539571
symmetry error     0.869297
fractal dimension error  0.379344
worst radius        0.649779
worst texture       0.386858
worst perimeter     0.674870
worst area          1.048970
worst smoothness    0.247199
worst compactness   0.915295
worst concavity     0.889174
worst concave points  0.492616
worst symmetry      0.521772
worst fractal dimension  0.831581
target             -0.528461
dtype: float64
```

since the data is still skewed, we have to apply transformation technique.

```
[25]: #to transform data using Square root transformation
df_transformed=df_capped.copy()
numericals=df_capped.select_dtypes(include=["number"])
for column in numericals.columns:
    df_transformed[column]=np.sqrt(numericals[column])
df_transformed.skew()
```

```
[25]: mean radius      0.423146
mean texture      0.180940
mean perimeter    0.456977
mean area         0.553411
mean smoothness   0.065613
mean compactness   0.372436
mean concavity    0.214155
mean concave points  0.039398
mean symmetry     0.391405
mean fractal dimension  0.566293
radius error       0.656397
texture error      0.297205
perimeter error    0.651761
area error         0.798276
smoothness error   0.415840
compactness error  0.463838
concavity error    -0.095279
concave points error  -0.822304
symmetry error     0.565571
fractal dimension error  0.536136
worst radius        0.612512
worst texture       0.109939
```

```
worst perimeter      0.628939
worst area          0.655672
worst smoothness    0.010433
worst compactness   0.344224
worst concavity    0.092076
worst concave points -0.443414
worst symmetry      0.243989
worst fractal dimension 0.635767
target              -0.528461
dtype: float64
```

```
[26]: corr_df_transformed.corr()
```

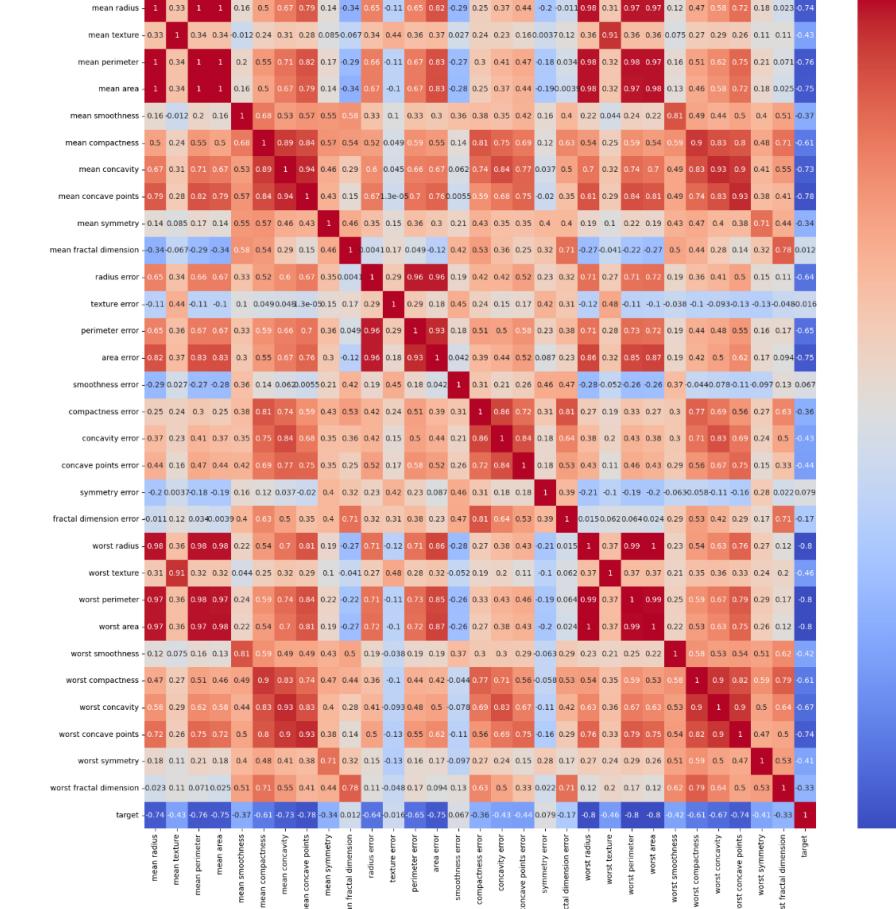
```
[26]:
```

	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	smo
mean radius	1.000000	0.331432	0.997845	0.997693	0.164113	0.499859	0.670608	0.791855	0.137927	-0.337797	...	0.308940	0.970047	0.971626	0.	
mean texture	0.331432	1.000000	0.337959	0.339155	-0.012197	0.242016	0.307556	0.276617	0.084836	-0.066660	...	0.914445	0.364612	0.360100	0.	
mean perimeter	0.997845	0.337959	1.000000	0.995357	0.201304	0.549985	0.708967	0.820942	0.171528	-0.286349	...	0.315611	0.976236	0.972685	0.	
mean area	0.997693	0.339155	0.995357	1.000000	0.159363	0.495928	0.671410	0.789275	0.140161	-0.335419	...	0.315553	0.971082	0.977411	0.	
mean smoothness	0.164113	-0.012197	0.201304	0.159363	1.000000	0.677998	0.528306	0.570484	0.552303	0.582067	...	0.044194	0.244414	0.216397	0.	
mean compactness	0.499859	0.242016	0.549985	0.495928	0.677998	1.000000	0.892266	0.840888	0.573652	0.542063	...	0.251987	0.594228	0.539349	0.	
mean concavity	0.670608	0.307556	0.708967	0.671410	0.528306	0.892266	1.000000	0.937028	0.456286	0.290317	...	0.319878	0.738102	0.699889	0.	
mean concave points	0.791855	0.276617	0.820942	0.789275	0.570484	0.840888	0.937028	1.000000	0.434842	0.152554	...	0.288815	0.836569	0.808324	0.	
mean symmetry	0.137927	0.084836	0.171528	0.140161	0.552303	0.573652	0.456286	0.434842	1.000000	0.461381	...	0.104028	0.215173	0.188338	0.	
mean fractal dimension	-0.337797	-0.066660	-0.286349	-0.335419	0.582067	0.542063	0.290317	0.152554	0.461381	1.000000	...	-0.040730	-0.222877	-0.268188	0.	
radius error	0.650093	0.341369	0.662758	0.668330	0.332057	0.524935	0.601269	0.670300	0.348775	0.004085	...	0.268726	0.706238	0.720982	0.	
texture error	-0.114931	0.442316	-0.105197	-0.103344	0.100393	0.048895	0.044585	-0.000013	0.147348	0.173819	...	0.481631	-0.107629	-0.104804	-0.	
perimeter error	0.654374	0.358651	0.673220	0.670477	0.331950	0.589063	0.656472	0.703249	0.355764	0.049297	...	0.283149	0.725480	0.715727	0.	
area error	0.817185	0.374078	0.825117	0.833179	0.296793	0.548491	0.670342	0.760327	0.299566	-0.115892	...	0.318443	0.854004	0.870812	0.	
smoothness error	-0.293128	0.026610	-0.272146	-0.277784	0.357722	0.139904	0.062222	0.005474	0.211585	0.416008	...	-0.052307	-0.262318	-0.259658	0.	
compactness error	0.252011	0.236295	0.299988	0.253473	0.379316	0.806996	0.738731	0.585190	0.432602	0.532937	...	0.187397	0.325628	0.270436	0.	
concavity error	0.373505	0.234160	0.412939	0.373747	0.349382	0.749162	0.844474	0.683889	0.353839	0.358901	...	0.204673	0.426419	0.379887	0.	
concave points error	0.439212	0.162450	0.470339	0.435049	0.423215	0.689283	0.766290	0.754564	0.347512	0.252411	...	0.113717	0.463326	0.425902	0.	
symmetry error	-0.199880	0.003730	-0.181274	-0.190674	0.155318	0.124085	0.037043	-0.020432	0.400831	0.319625	...	-0.104036	-0.191456	-0.201956	-0.	
fractal dimension error	-0.010952	0.122131	0.034178	-0.003890	0.395132	0.628846	0.500698	0.351937	0.400011	0.714518	...	0.061652	0.063868	0.024365	0.	
worst radius	0.975683	0.356884	0.976921	0.977434	0.220503	0.544115	0.700792	0.812718	0.185462	-0.272199	...	0.365005	0.994107	0.996804	0.	
worst texture	0.308940	0.914445	0.315611	0.315553	0.044194	0.251987	0.319878	0.288815	0.104028	-0.040730	...	1.000000	0.372883	0.365790	0.	
worst perimeter	0.970047	0.364612	0.976236	0.971082	0.244414	0.594228	0.738102	0.836569	0.215173	-0.222877	...	0.372883	1.000000	0.990409	0.	
worst area	0.971626	0.360100	0.972685	0.977411	0.216397	0.539349	0.699889	0.808324	0.188338	-0.268188	...	0.365790	0.990409	1.000000	0.	
worst smoothness	0.124388	0.074509	0.156180	0.126639	0.814667	0.588735	0.486732	0.488195	0.432076	0.499317	...	0.212161	0.245906	0.224990	1.	
worst compactness	0.469404	0.270758	0.513655	0.464452	0.493891	0.897525	0.832663	0.743051	0.465134	0.440353	...	0.347217	0.589572	0.530336	0.	
worst concavity	0.582224	0.286855	0.618254	0.579542	0.438550	0.830277	0.932962	0.829859	0.396127	0.275693	...	0.359639	0.672296	0.626435	0.	
worst concave points	0.724294	0.258143	0.752009	0.717043	0.496673	0.804585	0.900519	0.934142	0.382132	0.141389	...	0.329513	0.791873	0.753972	0.	
worst symmetry	0.181983	0.109775	0.206531	0.181845	0.396970	0.484258	0.405225	0.376211	0.714199	0.316293	...	0.241773	0.291882	0.263467	0.	
worst fractal dimension	0.023178	0.110732	0.070696	0.024748	0.505534	0.707800	0.547706	0.411244	0.444736	0.778991	...	0.200718	0.166815	0.117938	0.	
target	-0.740794	-0.425978	-0.755799	-0.751596	-0.368886	-0.609640	-0.727489	-0.775461	-0.335266	0.011787	...	-0.461307	-0.803972	-0.803976	-0.	

31 rows × 31 columns

```
[27]: plt.figure(figsize=(20,20))
sns.heatmap(corr, annot=True, cmap="coolwarm")
plt.show()
```

```
[27]:
```



some features are highly correlated each other such as worst radius, mean radius, mean perimeter, worst perimeter etc. most of the features have high positive correlation with the target.

```
[29]: from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
```

```
[30]: X=df_transformed.drop("target",axis=1)
X
```

```
[30]:
```

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	smot
-------------	--------------	----------------	-----------	-----------------	------------------	----------------	---------------------	---------------	------------------------	-----	--------------	---------------	-----------------	------------	------

Support Vector Machine (SVM)

How it Works:

- SVC uses a set of labeled training examples to find a decision boundary that separates the data points into different classes.
 - The decision boundary is represented as a linear function, and the goal is to find the boundary that maximizes the separation between the classes.

Why it's suitable:

- Works well when features are properly scaled.
 - Performs well for high-dimensional data, like this dataset with 30 features.
 - SVC is commonly used in image recognition, text classification, and medical diagnosis.

```
[55]: svc=SVC()
svc.fit(X_train,y_train)
```

```
[56]: svc_y_pred=svc.predict(X_test)  
svc_y_pred
```

Page 6

- ```

1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1.,
1., 1., 0., 1., 0., 1., 1., 0., 1., 1., 1., 1., 1., 1., 1., 1., 1., 0.,
0., 1., 1., 1., 1., 1., 1., 0., 0., 1., 1., 0., 0., 0., 1., 1., 1., 0., 0.,
1., 1., 0., 0., 1., 0., 1., 1., 1., 1., 1., 1., 1., 0., 1., 0., 0., 0., 0.,
0., 0., 0., 1., 1., 1., 1., 0., 1., 1., 1., 1., 1., 1., 0., 0., 1., 0., 0.,
0., 0., 0., 1., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 0., 1., 1., 0., 1.,
0., 0., 0., 1., 1., 1., 0., 1., 1., 0., 1., 1., 1., 1., 0., 1., 1., 1., 0.]

```

## k-Nearest Neighbors (k-NN)

## How it works:

- A non-parametric, instance-based algorithm that classifies a new point based on the majority class among its k closest neighbors.
  - Uses distance metrics (e.g., Euclidean distance) to determine proximity.

### Why it's suitable:

- it is Simple.
  - Can handle non-linear decision boundaries.
  - Sensitive to feature scaling and computationally expensive for large datasets

```
[59]: knn=KNeighborsClassifier(n_neighbors=5)
```

```
[59]: ▾ KNeighborsClassifier ⓘ ⓘ
KNeighborsClassifier()
```

```
[60]: knn_y_pred=knn.predict(x_test)
knn_y_pred
```

```
: array([1., 0., 0., 1., 1., 0., 0., 0.,
 1, 1, 0, 1, 1, 0, 1, 1])
```



### 3. Model Comparison (2 marks)

- Compare the performance of the five classification algorithms.
  - Which algorithm performed the best and which one performed the worst

```
[63]: from sklearn.metrics import accuracy_score,f1_score,precision_score,recall_score
```

```
[64]: predictions:[
 "Logistic Regression":logreg_y_pred,
 "Decision Tree Classifier":dtcl_y_pred,
 "Random Forest Classifier":rfcl_y_pred,
 "Support Vector Machine (SVM)":svc_y_pred,
 "k-Nearest Neighbors (k-NN)":knn_y_pred
}

result={
 "Model":[],
 "Accuracy Score":[],
 "F1 Score":[],
 "Precision Score":[],
 "Recall Score":[]
}

Compute metrics for each model

for model_name,y_pred in predictions.items():
 result["Model"].append(model_name)
 result["Accuracy Score"].append(accuracy_score(y_test,y_pred))
 result["F1 Score"].append(f1_score(y_test,y_pred))
 result["Precision Score"].append(precision_score(y_test,y_pred))
 result["Recall Score"].append(recall_score(y_test,y_pred))

print(result)

{'Model': ['Logistic Regression', 'Decision Tree Classifier', 'Random Forest Classifier', 'Support Vector Machine (SVM)', 'k-Nearest Neighbors (k-NN)'], 'Accuracy Score': [0.982456140598771, 0.947368421056315, 0.9649122807017544, 0.973684210563158, 0.956140350877193], 'F1 Score': [0.9859154929577465, 0.9722222222222222, 0.9790209790209791, 0.965934965034965], 'Precision Score': [0.9859154929577465, 0.97746478732394, 0.958804109580041, 0.9722222222222222, 0.9983333333333334], 'Recall Score': [0.9859154929577465, 0.957746478732394, 0.9859154929577465, 0.9859154929577465, 0.9718308591549313]}
```

| [65]: | Results                      |                |          |                 |              |
|-------|------------------------------|----------------|----------|-----------------|--------------|
| [65]: | Model                        | Accuracy Score | F1 Score | Precision Score | Recall Score |
| 0     | Logistic Regression          | 0.982456       | 0.985915 | 0.985915        | 0.985915     |
| 1     | Decision Tree Classifier     | 0.947368       | 0.957746 | 0.957746        | 0.957746     |
| 2     | Random Forest Classifier     | 0.964912       | 0.972222 | 0.958904        | 0.985915     |
| 3     | Support Vector Machine (SVM) | 0.973684       | 0.979021 | 0.972222        | 0.985915     |
| 4     | k-Nearest Neighbors (k-NN)   | 0.956140       | 0.965035 | 0.958333        | 0.971833     |

```
[66]: Results.sort_values(by="Accuracy Score", ascending=False)
```

| [66]: | Model                        | Accuracy Score | F1 Score | Precision Score | Recall Score |
|-------|------------------------------|----------------|----------|-----------------|--------------|
| 0     | Logistic Regression          | 0.982456       | 0.985915 | 0.985915        | 0.985915     |
| 3     | Support Vector Machine (SVM) | 0.973684       | 0.979021 | 0.972222        | 0.985915     |
| 2     | Random Forest Classifier     | 0.964912       | 0.972222 | 0.958904        | 0.985915     |
| 4     | k-Nearest Neighbors (k-NN)   | 0.956140       | 0.965035 | 0.958333        | 0.971333     |
| 1     | Decision Tree Classifier     | 0.947368       | 0.957746 | 0.957746        | 0.957746     |

- From the above result table, Logistic Regression is the best model with highest accuracy score, F1 score, precision score and recall score.
  - The worst model is Decision Tree Classifier which has the lowest accuracy score compared to other models.

上一頁 下一頁