

```
In [2]: # import the data file
data = pd.read_csv('data.csv')

In [3]: # display the first 5 data
data.head()

Out[3]:
   id  diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  smoothness_mean  compactness_mean  concavity_mean  concave
points_mean  ...  texture_worst  perimeter_worst
0   842302      M         17.99         10.38         122.80        1001.0         0.11840         0.27760         0.3001         0.14710
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
1   842517      M         20.57         17.77         132.90        1326.0         0.08474         0.07864         0.0869         0.07017
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
2   8430903     M         19.69         21.25         130.00        1203.0         0.10960         0.15990         0.1974         0.12790
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
3   8434301     M         11.42         20.38         77.58         386.1         0.14250         0.28390         0.2414         0.10520
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
4   8435802     M         20.29         14.34         135.10        1297.0         0.10030         0.13280         0.1980         0.10430
...         ...         ...         ...         ...         ...         ...         ...         ...         ...

5 rows × 33 columns

In [4]: import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

In [5]: # display the last 5 data
data.tail()

Out[5]:
   id  diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  smoothness_mean  compactness_mean  concavity_mean  concave
points_mean  ...  texture_worst  perimeter_worst
564  926424      M         21.56         22.39         142.00        1479.0         0.11100         0.11590         0.24390         0.13890
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
565  926682      M         20.13         28.25         131.20        1261.0         0.09780         0.10340         0.14400         0.09781
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
566  926954      M         16.60         28.08         108.30         858.1         0.08455         0.10230         0.09251         0.05302
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
567  927241      M         20.60         29.33         140.10        1265.0         0.11780         0.27700         0.35140         0.15200
...         ...         ...         ...         ...         ...         ...         ...         ...         ...
568  92751      B          7.76         24.54         47.92         181.0         0.05263         0.04362         0.00000         0.00000
...         ...         ...         ...         ...         ...         ...         ...         ...         ...

5 rows × 33 columns

In [6]: # count the number of rows and columns in the dataset
data.shape

Out[6]: (569, 33)

In [7]: # count the number of empty [NaN, NAN , na] values in each columns
data.isna().sum()

Out[7]:
id                    0
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
unnamed: 32          569
dtype: int64

In [8]: # Drop the column with all missing values [Input can be 0 or 1 for Integer and 'index' or 'columns' for String]
data = data.dropna(axis=1)

In [9]: # Get the new count all the number of rows and cols again
data.shape

Out[9]: (569, 32)

In [10]: # Get a count of the number of Malignant(M) or Benign(B) cells
data['diagnosis'].value_counts()

Out[10]:
B    357
M    212
Name: diagnosis, dtype: int64

In [11]: # Visualize the count of M and B in the dataset x => data['diagnosis'], y => label='count'
sns.countplot(data['diagnosis'], label='count')

C:\Users\taluk\anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version
0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misint
erpretation.
  warnings.warn(

Out[11]:
<AxesSubplot: xlabel='diagnosis', ylabel='count'>

```