

Problem Statement:

Breast Cancer Prediction

In [1]:

```
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
```

Data collection:

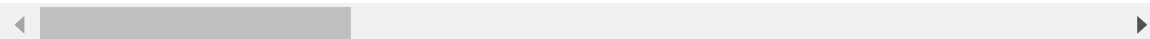
In [2]:

```
df=pd.read_csv(r"C:\Users\raja\Downloads\BreastCancerPrediction.csv")
df
```

Out[2]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothn
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	
...	
564	926424	M	21.56	22.39	142.00	1479.0	
565	926682	M	20.13	28.25	131.20	1261.0	
566	926954	M	16.60	28.08	108.30	858.1	
567	927241	M	20.60	29.33	140.10	1265.0	
568	92751	B	7.76	24.54	47.92	181.0	

569 rows × 33 columns



Data Cleaning and Preprocessing

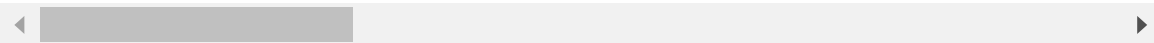
In [3]:

```
df.head()
```

Out[3]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
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	

5 rows × 33 columns



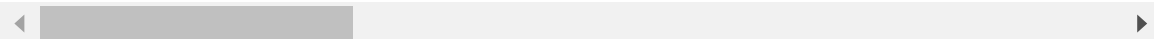
In [4]:

```
df.tail()
```

Out[4]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
564	926424	M	21.56	22.39	142.00	1479.0	
565	926682	M	20.13	28.25	131.20	1261.0	
566	926954	M	16.60	28.08	108.30	858.1	
567	927241	M	20.60	29.33	140.10	1265.0	
568	92751	B	7.76	24.54	47.92	181.0	

5 rows × 33 columns



In [5]:

```
df.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
```

In [6]:

```
df.isnull().sum()
```

Out[6]:

```
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
```

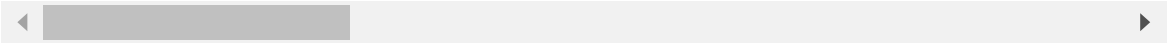
In [7]:

```
df.drop(['Unnamed: 32'],axis=1)
```

Out[7]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothn
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	
...	
564	926424	M	21.56	22.39	142.00	1479.0	
565	926682	M	20.13	28.25	131.20	1261.0	
566	926954	M	16.60	28.08	108.30	858.1	
567	927241	M	20.60	29.33	140.10	1265.0	
568	92751	B	7.76	24.54	47.92	181.0	

569 rows × 32 columns

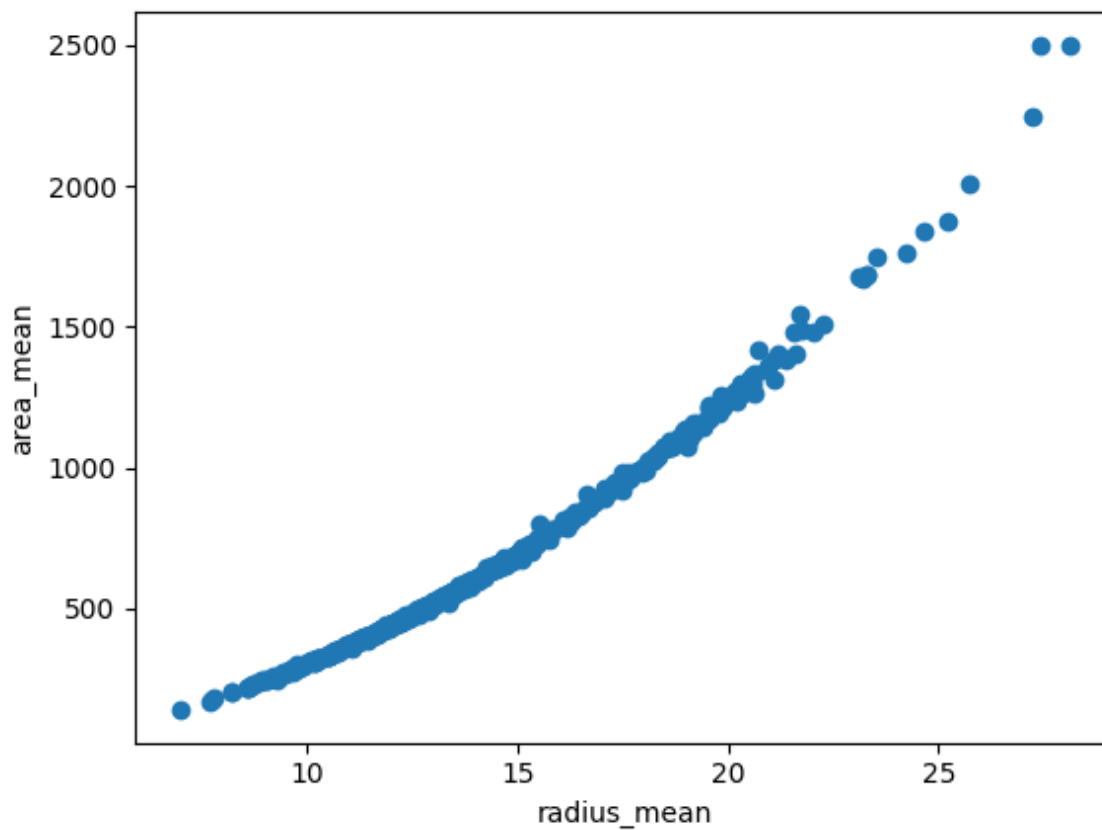


In [8]:

```
plt.scatter(df["radius_mean"],df["area_mean"])  
plt.xlabel("radius_mean")  
plt.ylabel("area_mean")
```

Out[8]:

Text(0, 0.5, 'area_mean')



KMeans Clustering

In [9]:

```
from sklearn.cluster import KMeans
```

In [10]:

```
km=KMeans()  
km
```

Out[10]:

```
▼ KMeans  
KMeans()
```

In [11]:

```
y_predicted=km.fit_predict(df[["radius_mean","area_mean"]])
y_predicted
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

Out[11]:

```
array([4, 1, 1, 5, 1, 5, 4, 0, 0, 5, 6, 6, 1, 6, 0, 6, 6, 6, 1, 0, 0, 2,
        6, 1, 4, 4, 6, 4, 6, 4, 4, 5, 4, 1, 6, 4, 0, 0, 6, 0, 0, 5, 4, 0,
        0, 4, 2, 0, 5, 0, 5, 0, 5, 4, 6, 5, 1, 6, 0, 2, 2, 2, 6, 2, 0, 6,
        2, 5, 2, 0, 1, 2, 4, 0, 5, 6, 0, 4, 1, 0, 5, 0, 7, 1, 5, 4, 6, 4,
        5, 6, 6, 6, 0, 0, 6, 1, 5, 2, 5, 6, 0, 2, 5, 2, 2, 0, 5, 5, 7, 5,
        2, 5, 0, 2, 2, 5, 2, 6, 6, 4, 5, 4, 7, 6, 0, 0, 0, 1, 6, 1, 5, 6,
        6, 6, 4, 0, 5, 5, 6, 5, 2, 6, 5, 0, 5, 5, 5, 6, 6, 0, 0, 2, 2, 5,
        0, 5, 4, 4, 5, 5, 5, 1, 1, 5, 7, 6, 5, 4, 4, 6, 5, 0, 6, 5, 2, 2,
        2, 6, 0, 0, 3, 1, 6, 5, 6, 2, 4, 5, 5, 5, 0, 0, 2, 5, 6, 0, 0, 4,
        1, 6, 5, 4, 7, 0, 5, 6, 2, 4, 0, 6, 1, 5, 3, 4, 0, 0, 5, 2, 1, 1,
        0, 0, 2, 6, 0, 6, 2, 6, 0, 0, 4, 5, 5, 1, 2, 0, 7, 1, 0, 4, 0, 5,
        5, 0, 1, 2, 0, 0, 2, 5, 1, 5, 1, 4, 1, 0, 1, 6, 6, 6, 1, 4, 4, 6,
        4, 1, 2, 0, 0, 2, 0, 5, 7, 2, 4, 5, 5, 4, 0, 0, 1, 5, 1, 6, 0, 0,
        5, 0, 5, 5, 6, 6, 0, 5, 0, 0, 5, 5, 0, 2, 1, 5, 1, 2, 5, 5, 0, 2,
        0, 0, 5, 6, 0, 5, 2, 5, 5, 4, 2, 5, 2, 1, 0, 1, 5, 0, 0, 5, 6, 6,
        6, 0, 5, 5, 5, 4, 0, 4, 2, 7, 6, 2, 5, 1, 5, 2, 5, 6, 5, 5, 5, 6,
        3, 6, 5, 5, 0, 0, 2, 2, 0, 0, 0, 6, 0, 1, 1, 5, 7, 7, 6, 6, 1, 1,
        0, 6, 2, 0, 0, 5, 5, 5, 5, 5, 0, 6, 5, 0, 5, 1, 2, 2, 6, 1, 5, 0,
        0, 0, 5, 5, 4, 5, 0, 0, 5, 5, 6, 0, 4, 5, 5, 5, 2, 6, 6, 5, 2, 6,
        0, 5, 5, 6, 5, 0, 2, 2, 2, 5, 5, 0, 6, 5, 1, 4, 6, 0, 0, 0, 0, 0,
        5, 4, 0, 2, 4, 5, 4, 6, 6, 1, 5, 1, 5, 6, 0, 0, 5, 0, 0, 2, 4, 3,
        6, 5, 0, 0, 0, 2, 4, 5, 2, 5, 6, 5, 5, 0, 0, 0, 5, 6, 5, 0, 0, 0,
        6, 5, 6, 1, 5, 4, 5, 4, 4, 5, 0, 6, 5, 5, 4, 1, 6, 0, 5, 7, 2, 2,
        5, 5, 6, 6, 5, 6, 0, 6, 6, 5, 4, 1, 0, 0, 2, 7, 5, 0, 2, 2, 0, 5,
        0, 5, 5, 5, 0, 1, 5, 1, 0, 5, 2, 2, 5, 6, 6, 0, 0, 0, 2, 2, 2, 5,
        5, 5, 0, 2, 0, 2, 2, 2, 6, 5, 0, 5, 6, 1, 7, 1, 4, 1, 2])
```

In [12]:

```
df["Cluster"]=y_predicted
df.head()
```

Out[12]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
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	

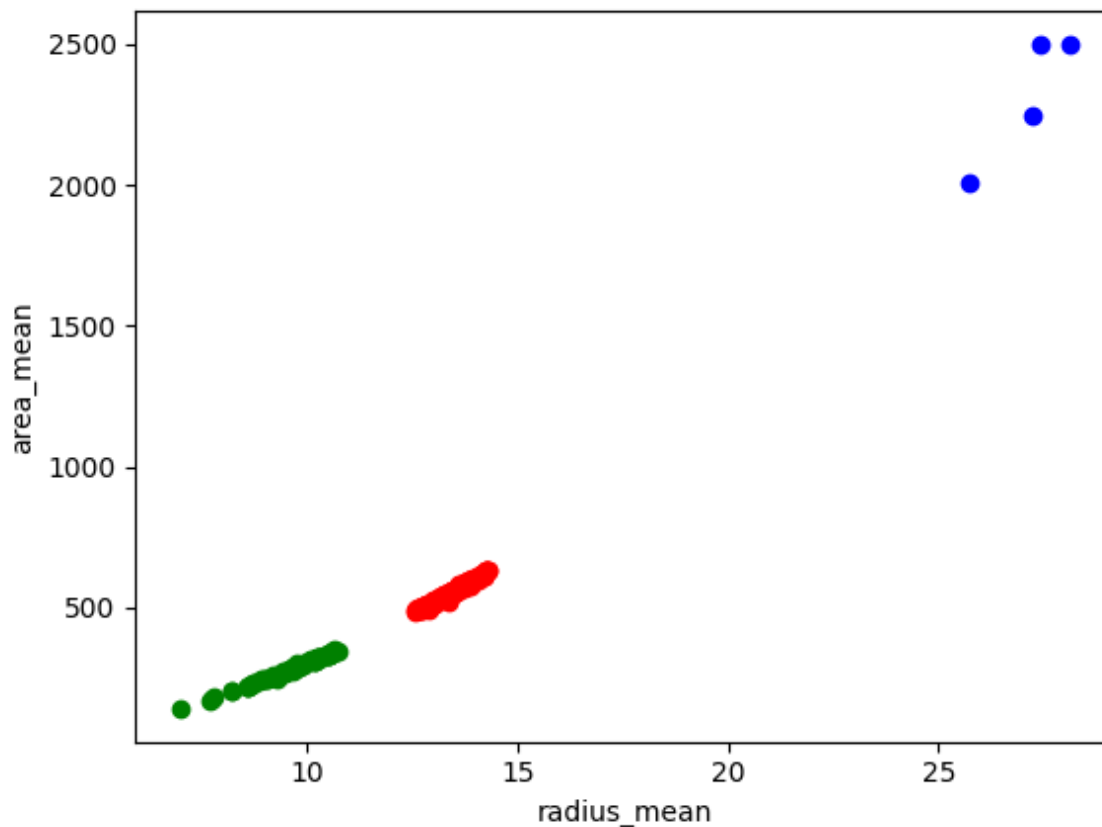
5 rows × 34 columns

In [13]:

```
df1=df[df.Cluster==0]
df2=df[df.Cluster==2]
df3=df[df.Cluster==3]
plt.scatter(df1["radius_mean"],df1["area_mean"],color="red")
plt.scatter(df2["radius_mean"],df2["area_mean"],color="green")
plt.scatter(df3["radius_mean"],df3["area_mean"],color="blue")
plt.xlabel("radius_mean")
plt.ylabel("area_mean")
```

Out[13]:

Text(0, 0.5, 'area_mean')



In [14]:

```
from sklearn.preprocessing import MinMaxScaler
```

In [15]:

```
scaler=MinMaxScaler()
```


In [16]:

```
scaler.fit(df[["area_mean"]])  
df["area_mean"]=scaler.transform(df[["area_mean"]])  
df.head()
```

Out[16]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	17.99	10.38	122.80	0.363733	
1	842517	M	20.57	17.77	132.90	0.501591	
2	84300903	M	19.69	21.25	130.00	0.449417	
3	84348301	M	11.42	20.38	77.58	0.102906	
4	84358402	M	20.29	14.34	135.10	0.489290	

5 rows × 34 columns



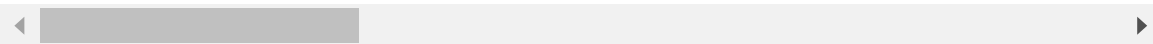
In [17]:

```
scaler.fit(df[["radius_mean"]])  
df["radius_mean"]=scaler.transform(df[["radius_mean"]])  
df.head()
```

Out[17]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	0.521037	10.38	122.80	0.363733	
1	842517	M	0.643144	17.77	132.90	0.501591	
2	84300903	M	0.601496	21.25	130.00	0.449417	
3	84348301	M	0.210090	20.38	77.58	0.102906	
4	84358402	M	0.629893	14.34	135.10	0.489290	

5 rows × 34 columns



In [18]:

```
km=KMeans()
```

In [19]:

```
y_predicted=km.fit_predict(df[["radius_mean","area_mean"]])
y_predicted
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

Out[19]:

```
array([2, 3, 3, 1, 3, 1, 2, 7, 7, 1, 5, 5, 3, 5, 7, 5, 5, 5, 3, 7, 7, 4,
       5, 3, 2, 2, 5, 2, 5, 2, 2, 1, 2, 3, 5, 2, 7, 7, 5, 7, 7, 1, 3, 7,
       7, 2, 4, 7, 1, 7, 1, 7, 1, 2, 5, 1, 3, 5, 7, 4, 4, 4, 5, 4, 7, 5,
       4, 1, 4, 7, 3, 4, 2, 7, 1, 5, 7, 2, 3, 7, 1, 7, 0, 3, 1, 2, 5, 2,
       1, 5, 5, 5, 7, 7, 5, 3, 1, 4, 1, 5, 7, 4, 1, 4, 4, 7, 1, 1, 0, 1,
       4, 7, 7, 4, 4, 1, 4, 5, 5, 2, 1, 2, 0, 5, 7, 7, 7, 3, 5, 3, 1, 5,
       5, 5, 2, 7, 1, 1, 5, 1, 4, 5, 1, 7, 1, 1, 1, 5, 5, 7, 7, 4, 4, 1,
       7, 1, 2, 2, 1, 1, 1, 3, 3, 1, 0, 5, 1, 2, 2, 5, 1, 7, 5, 1, 1, 4,
       4, 5, 7, 7, 6, 3, 5, 1, 5, 4, 2, 1, 1, 1, 7, 7, 4, 1, 5, 7, 7, 2,
       3, 5, 1, 2, 0, 7, 1, 5, 4, 2, 7, 5, 3, 1, 6, 2, 7, 7, 1, 4, 3, 3,
       7, 7, 4, 5, 7, 5, 4, 5, 7, 7, 2, 1, 1, 3, 4, 7, 0, 3, 7, 2, 7, 1,
       1, 7, 3, 4, 7, 7, 1, 1, 3, 1, 3, 2, 3, 7, 3, 5, 5, 5, 3, 2, 2, 5,
       2, 3, 4, 7, 7, 1, 7, 1, 0, 4, 2, 1, 1, 2, 7, 7, 3, 1, 3, 5, 7, 7,
       1, 7, 1, 1, 5, 5, 7, 1, 7, 7, 1, 1, 7, 4, 3, 1, 3, 4, 1, 1, 7, 4,
       7, 7, 1, 5, 7, 1, 4, 1, 1, 2, 4, 1, 4, 3, 7, 3, 1, 7, 7, 1, 5, 5,
       5, 7, 1, 1, 1, 2, 7, 2, 4, 0, 5, 4, 1, 3, 1, 4, 1, 5, 1, 1, 1, 5,
       6, 5, 1, 7, 7, 7, 4, 4, 7, 7, 7, 5, 7, 3, 3, 1, 0, 0, 5, 5, 3, 3,
       7, 5, 4, 7, 7, 1, 1, 1, 1, 1, 7, 5, 1, 7, 1, 3, 4, 4, 5, 3, 1, 7,
       7, 7, 1, 1, 2, 1, 7, 7, 1, 1, 5, 7, 2, 1, 1, 1, 4, 5, 5, 1, 4, 5,
       7, 1, 1, 5, 1, 7, 4, 4, 4, 1, 1, 7, 5, 1, 3, 2, 5, 7, 7, 7, 7, 7,
       1, 2, 7, 4, 2, 1, 2, 5, 5, 3, 1, 3, 1, 5, 7, 7, 1, 7, 7, 4, 2, 6,
       5, 1, 7, 7, 7, 4, 2, 1, 4, 1, 5, 1, 1, 7, 7, 7, 1, 5, 1, 7, 7, 7,
       5, 1, 5, 3, 1, 2, 1, 2, 2, 1, 7, 5, 7, 1, 2, 3, 5, 7, 1, 0, 4, 4,
       1, 1, 5, 5, 1, 5, 7, 5, 5, 1, 2, 3, 7, 7, 4, 0, 1, 7, 4, 4, 7, 1,
       7, 1, 1, 1, 7, 3, 1, 3, 7, 1, 4, 4, 1, 5, 5, 7, 7, 7, 4, 4, 1,
       1, 1, 7, 4, 7, 4, 4, 4, 5, 1, 7, 1, 5, 3, 0, 3, 2, 3, 4])
```

In [20]:

```
df["New cluster"]=y_predicted
df.head()
```

Out[20]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	0.521037	10.38	122.80	0.363733	
1	842517	M	0.643144	17.77	132.90	0.501591	
2	84300903	M	0.601496	21.25	130.00	0.449417	
3	84348301	M	0.210090	20.38	77.58	0.102906	
4	84358402	M	0.629893	14.34	135.10	0.489290	

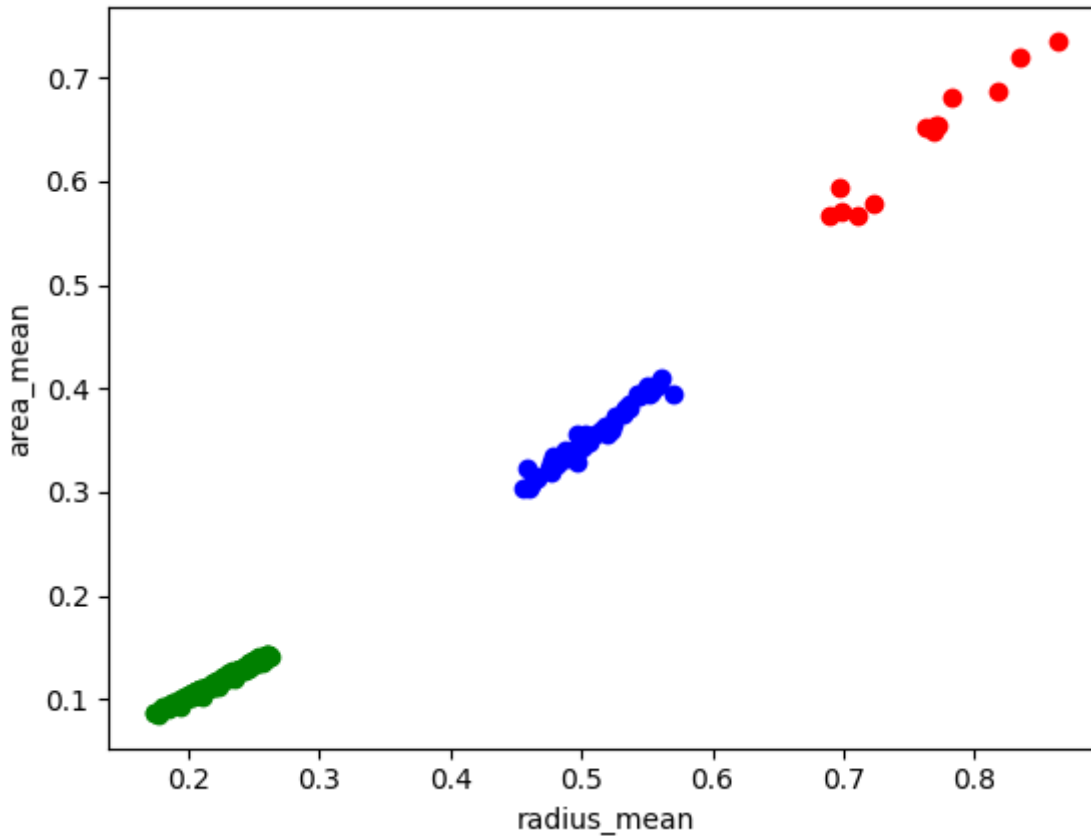
5 rows × 35 columns

In [21]:

```
df1=df[df["New cluster"]==0]
df2=df[df["New cluster"]==1]
df3=df[df["New cluster"]==2]
plt.scatter(df1["radius_mean"],df1["area_mean"],color="red")
plt.scatter(df2["radius_mean"],df2["area_mean"],color="green")
plt.scatter(df3["radius_mean"],df3["area_mean"],color="blue")
plt.xlabel("radius_mean")
plt.ylabel("area_mean")
```

Out[21]:

Text(0, 0.5, 'area_mean')



In [22]:

```
km.cluster_centers_
```

Out[22]:

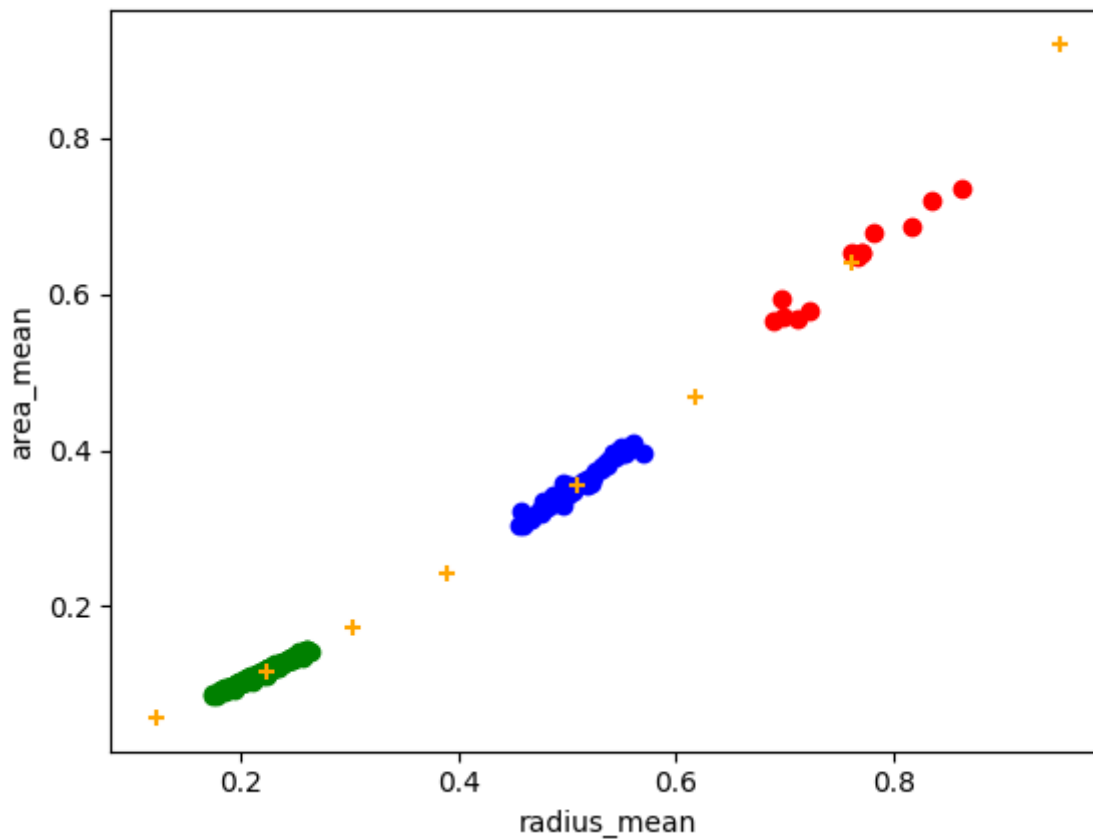
```
array([[0.7609556 , 0.63925279],
       [0.2231628 , 0.11753918],
       [0.50987888, 0.35605759],
       [0.61719997, 0.46705485],
       [0.12171243, 0.05773363],
       [0.38962757, 0.243164  ],
       [0.95314497, 0.92110286],
       [0.30182381, 0.17259778]])
```

In [23]:

```
df1=df[df["New cluster"]==0]
df2=df[df["New cluster"]==1]
df3=df[df["New cluster"]==2]
plt.scatter(df1["radius_mean"],df1["area_mean"],color="red")
plt.scatter(df2["radius_mean"],df2["area_mean"],color="green")
plt.scatter(df3["radius_mean"],df3["area_mean"],color="blue")
plt.scatter(km.cluster_centers_[0],km.cluster_centers_[1],color="orange",marker="+")
plt.xlabel("radius_mean")
plt.ylabel("area_mean")
```

Out[23]:

Text(0, 0.5, 'area_mean')



In [24]:

```
k_rng=range(1,10)
sse=[]
```

In [25]:

```

for k in k_rng:
    km=KMeans(n_clusters=k)
    km.fit(df[["radius_mean", "area_mean"]])
    sse.append(km.inertia_)
#km.inertia_ will give you the value of sum of square errorprint(sse)
plt.plot(k_rng,sse)
plt.xlabel("K")
plt.ylabel("Sum of Squared Error")

```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

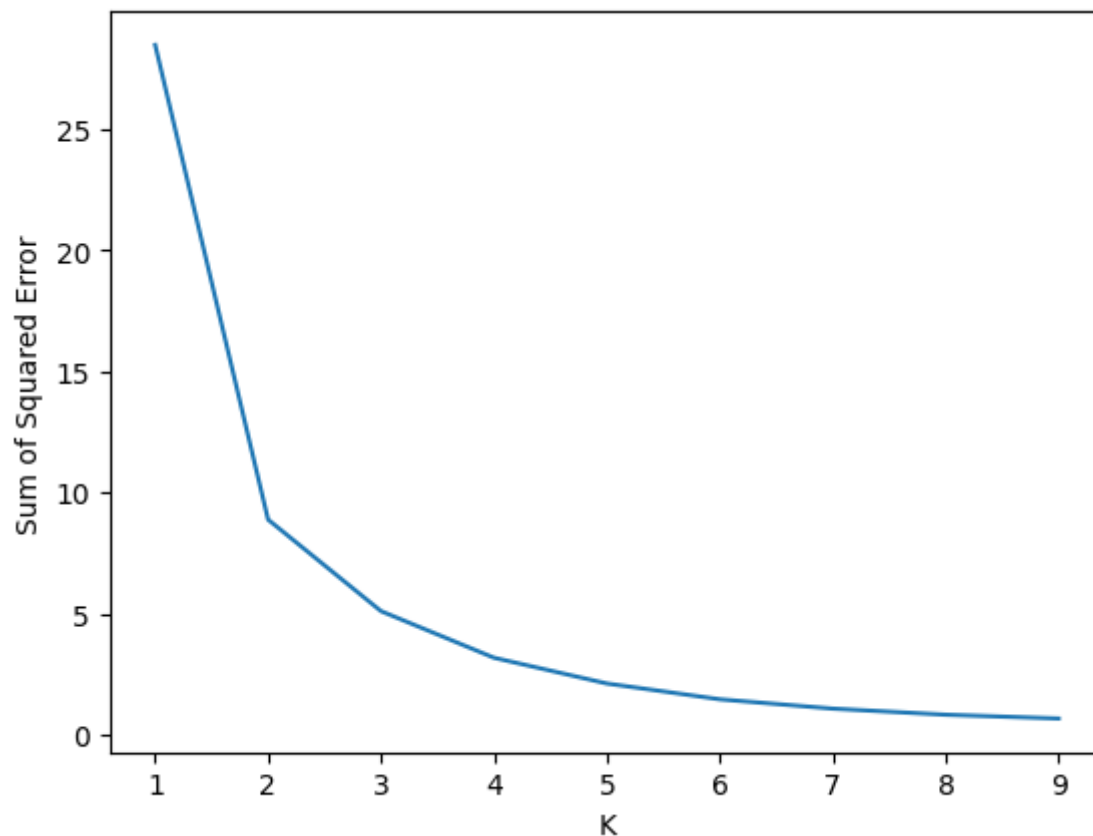
```
warnings.warn(
```

C:\Users\raja\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\cluster_kmeans.py:870: FutureWarning: The default value of `n_init` will change from 10 to 'auto' in 1.4. Set the value of `n_init` explicitly to suppress the warning

```
warnings.warn(
```

Out[25]:

```
Text(0, 0.5, 'Sum of Squared Error')
```



Conclusion:

In this dataset we are doing clustering on Radius_mean and Area_mean. This is the best model for this dataset. When k value is high error rate is low, or k value is low error rate is high.