

Problem statement:

Predictive study using the breast cancer diagnosis dataset.

1.Data Collection

In [43]:

```
import pandas as pd
from matplotlib import pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
```

In [44]:

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

Out[44]:

	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



In [45]:

```
df.shape
```

Out[45]:

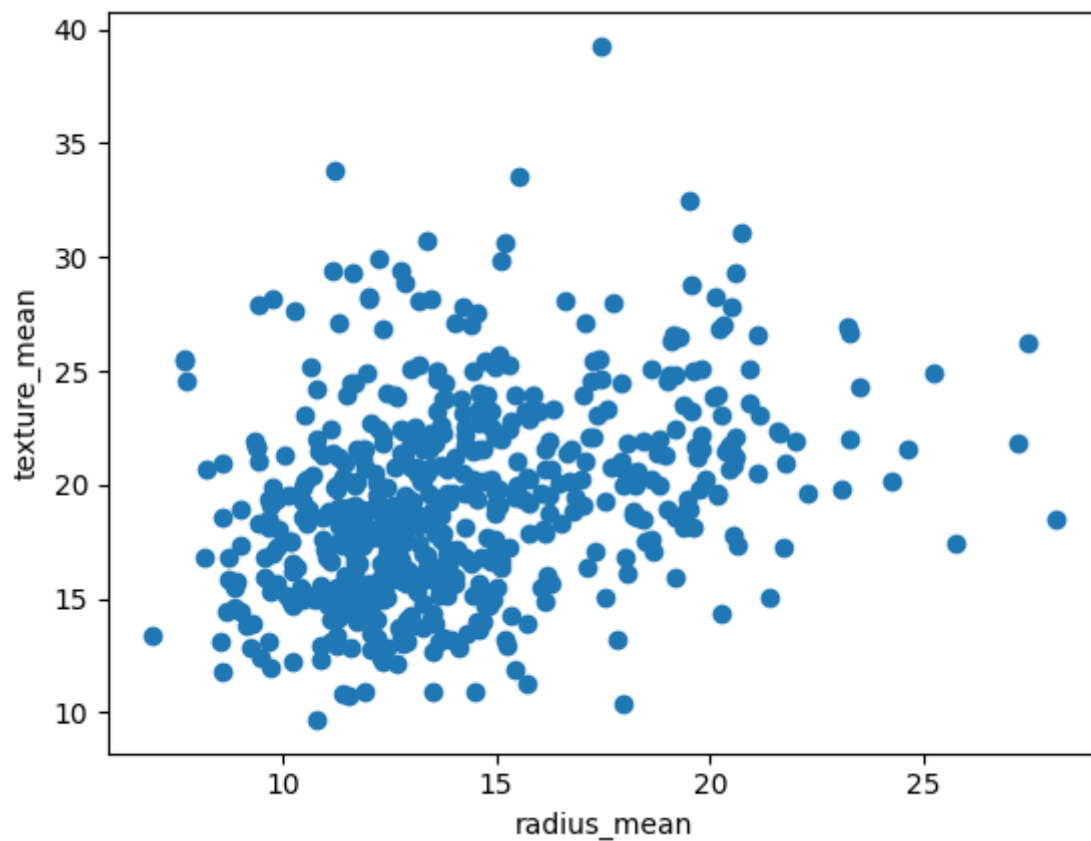
(569, 33)

In [46]:

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

Out[46]:

Text(0, 0.5, 'texture_mean')



In [47]:

```
from sklearn.cluster import KMeans  
km=KMeans()  
km
```

Out[47]:

▼ KMeans
KMeans()

In [48]:

```
y_pred=km.fit_predict(df[["radius_mean","texture_mean"]])
y_pred
```

Out[48]:

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

In [49]:

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

Out[49]:

	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

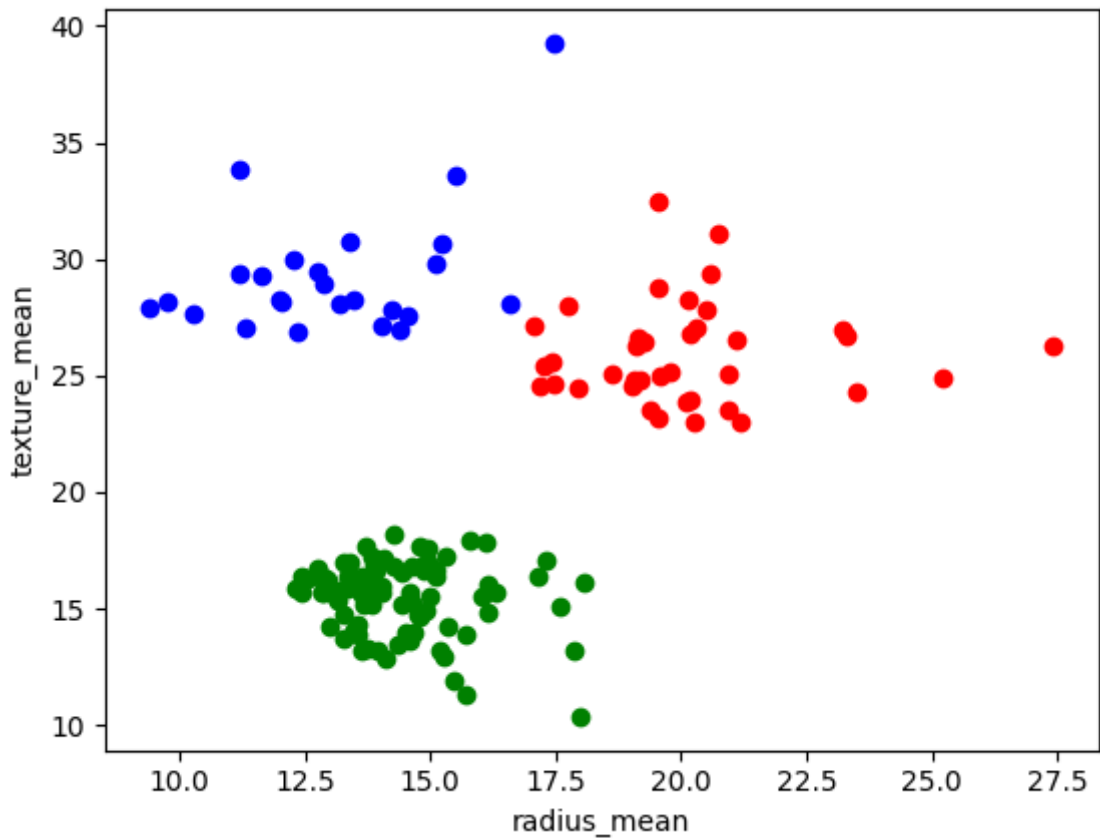


In [50]:

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

Out[50]:

Text(0, 0.5, 'texture_mean')



In [51]:

```
from sklearn.preprocessing import MinMaxScaler
scaler=MinMaxScaler()
scaler.fit(df[["texture_mean"]])
df["texture_mean"]=scaler.transform(df[["texture_mean"]])
df.head()
```

Out[51]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	17.99	0.022658	122.80	1001.0	
1	842517	M	20.57	0.272574	132.90	1326.0	
2	84300903	M	19.69	0.390260	130.00	1203.0	
3	84348301	M	11.42	0.360839	77.58	386.1	
4	84358402	M	20.29	0.156578	135.10	1297.0	

5 rows × 34 columns



In [52]:

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

Out[52]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	0.521037	0.022658	122.80	1001.0	
1	842517	M	0.643144	0.272574	132.90	1326.0	
2	84300903	M	0.601496	0.390260	130.00	1203.0	
3	84348301	M	0.210090	0.360839	77.58	386.1	
4	84358402	M	0.629893	0.156578	135.10	1297.0	

5 rows × 34 columns



In [53]:

```
y_pred=km.fit_predict(df[["radius_mean","texture_mean"]])
y_pred
```

Out[53]:

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

In [54]:

```
df["New Cluster"]=y_pred
df.head()
```

Out[54]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothnes
0	842302	M	0.521037	0.022658	122.80	1001.0	
1	842517	M	0.643144	0.272574	132.90	1326.0	
2	84300903	M	0.601496	0.390260	130.00	1203.0	
3	84348301	M	0.210090	0.360839	77.58	386.1	
4	84358402	M	0.629893	0.156578	135.10	1297.0	

5 rows × 35 columns

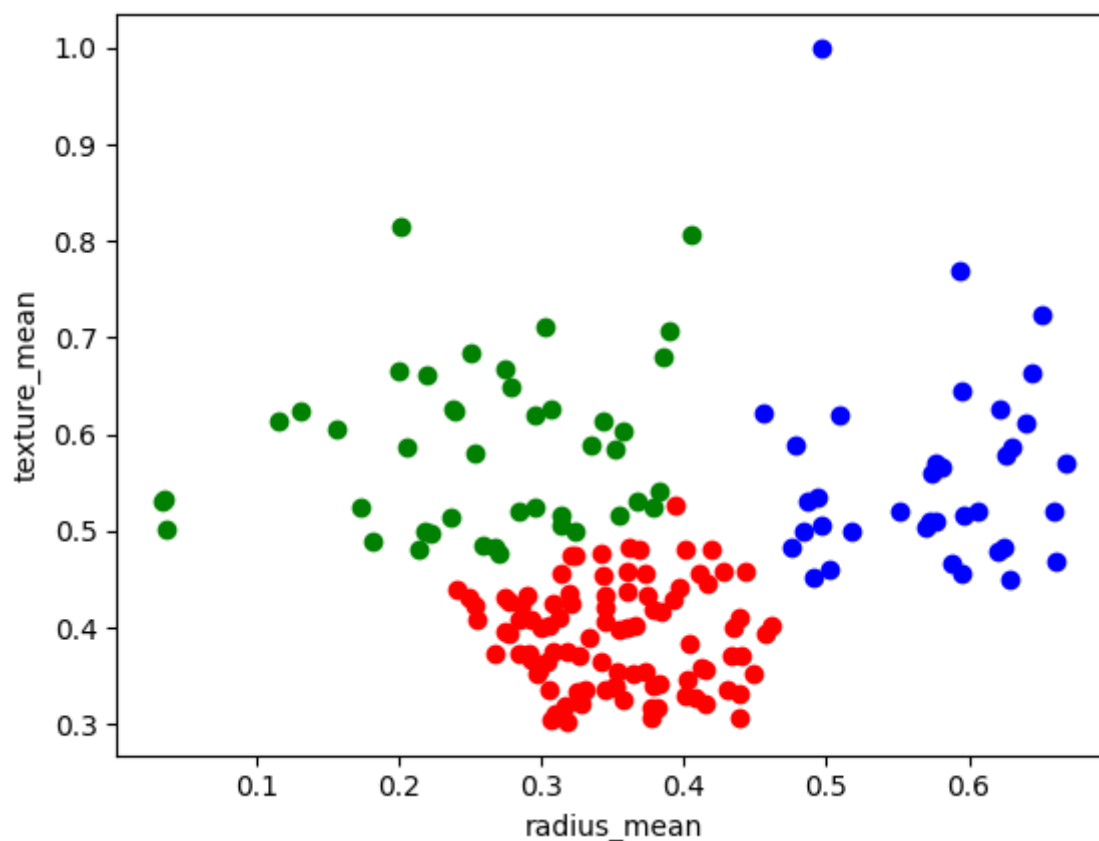


In [55]:

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

Out[55]:

Text(0, 0.5, 'texture_mean')



In [56]:

```
km.cluster_centers_
```

Out[56]:

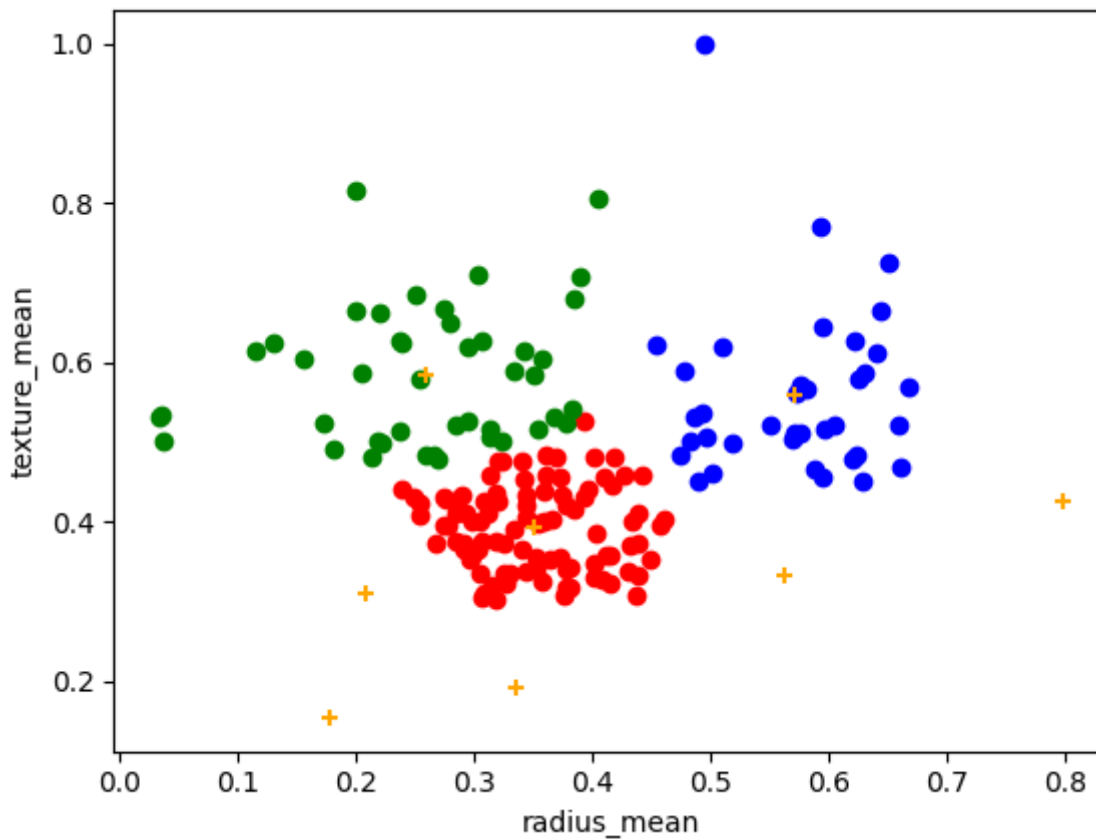
```
array([[0.35173159, 0.39188367],
       [0.79840767, 0.42469846],
       [0.2590623 , 0.58293879],
       [0.57132058, 0.55893025],
       [0.33570532, 0.19063107],
       [0.20867092, 0.3094643 ],
       [0.17750575, 0.15412045],
       [0.56287997, 0.33184226]])
```

In [57]:

```
df1=df[df["New Cluster"]==0]
df2=df[df["New Cluster"]==2]
df3=df[df["New Cluster"]==3]
plt.scatter(df1["radius_mean"],df1["texture_mean"],color="red")
plt.scatter(df2["radius_mean"],df2["texture_mean"],color="green")
plt.scatter(df3["radius_mean"],df3["texture_mean"],color="blue")
plt.scatter(km.cluster_centers_[0],km.cluster_centers_[1],color="orange",marker="+")
plt.xlabel("radius_mean")
plt.ylabel("texture_mean")
```

Out[57]:

Text(0, 0.5, 'texture_mean')



In [58]:

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

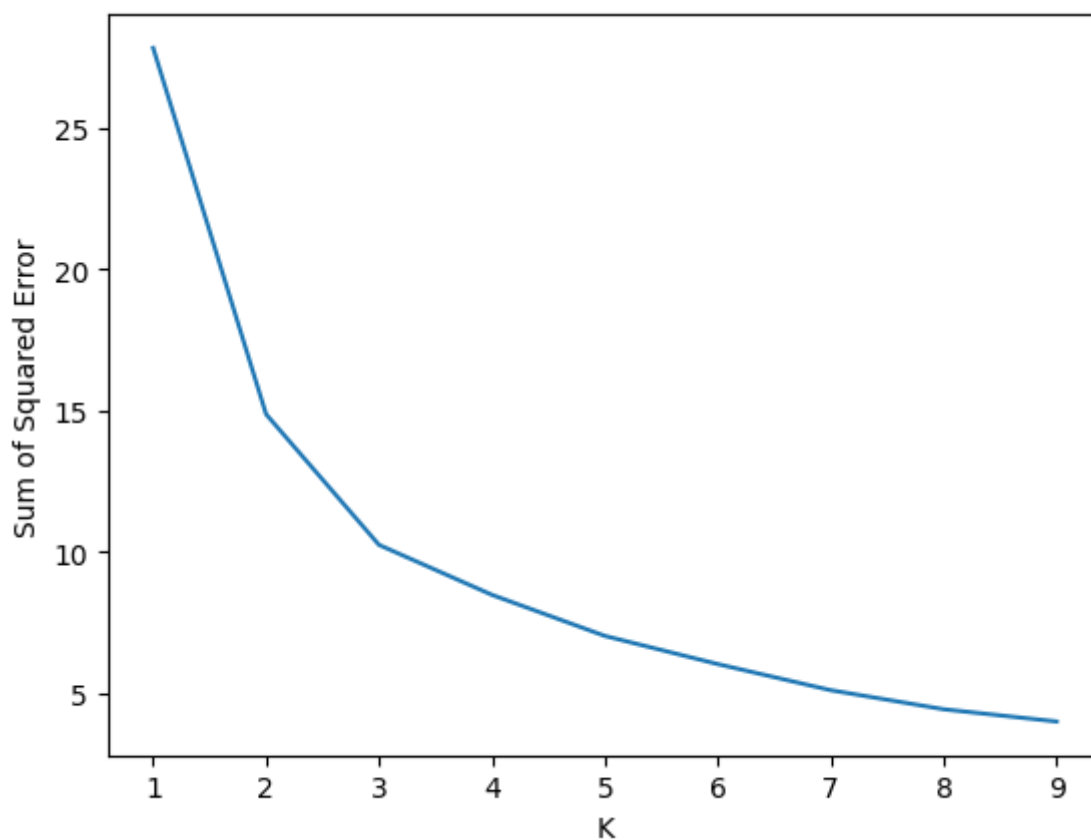

In [59]:

```
for k in k_rng:
    km=KMeans(n_clusters=k)
    km.fit(df[["radius_mean", "texture_mean"]])
    sse.append(km.inertia_)
print(sse)
plt.plot(k_rng,sse)
plt.xlabel("K")
plt.ylabel("Sum of Squared Error")
```

```
[27.81750759504307, 14.87203295827117, 10.252751496105196, 8.4876127848450
78, 7.035500433198194, 6.039305768835716, 5.116755795030002, 4.44301570025
843, 4.010136919135096]
```

Out[59]:

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



Conclusion:

For the given dataset we can do prediction by various models, but accuracy from those models is not good. So we prefer K-Means Clustering for this dataset.