

Problem statement:-

Predictive study using the breast cancer diagnosis dataset

Data Collection

In [1]:

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

In [2]:

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

Out[2]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_r
0	842302	M	17.99	10.38	122.80	1001.0	0.1
1	842517	M	20.57	17.77	132.90	1326.0	0.0
2	84300903	M	19.69	21.25	130.00	1203.0	0.1
3	84348301	M	11.42	20.38	77.58	386.1	0.1
4	84358402	M	20.29	14.34	135.10	1297.0	0.1
...
564	926424	M	21.56	22.39	142.00	1479.0	0.1
565	926682	M	20.13	28.25	131.20	1261.0	0.0
566	926954	M	16.60	28.08	108.30	858.1	0.0
567	927241	M	20.60	29.33	140.10	1265.0	0.1
568	92751	B	7.76	24.54	47.92	181.0	0.0

569 rows × 33 columns



In [3]:

```
df.shape
```

Out[3]:

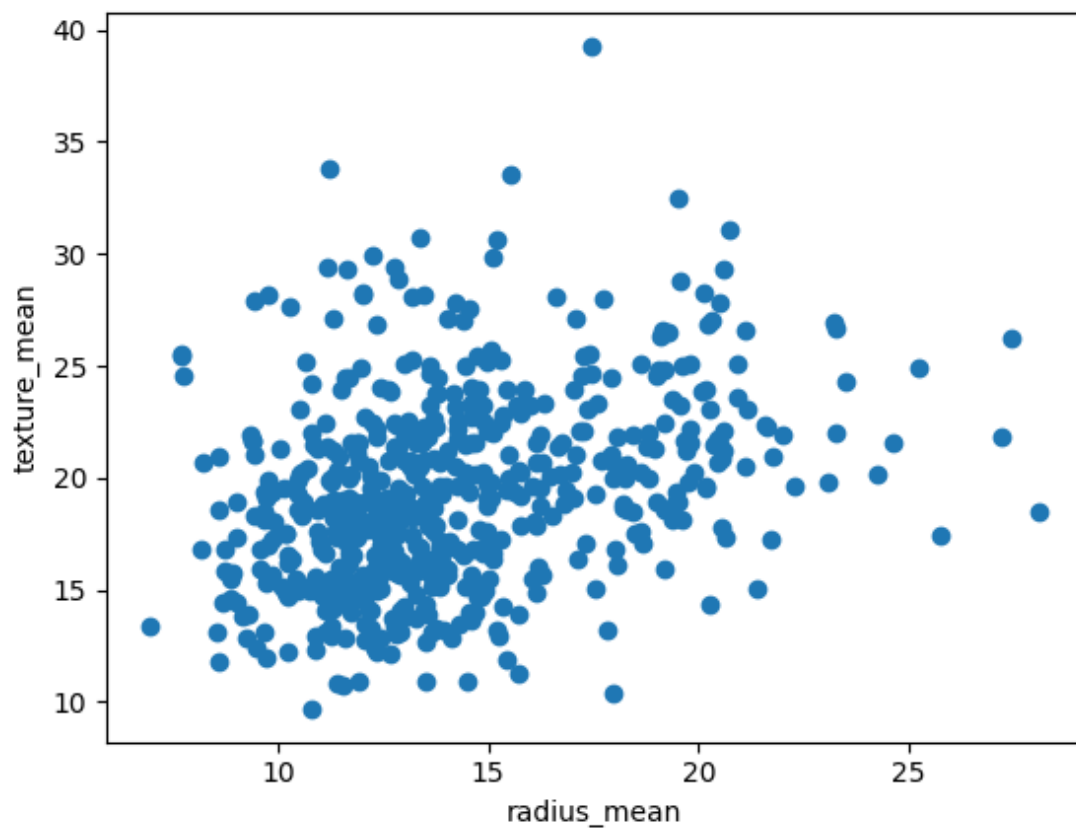
(569, 33)

In [4]:

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

Out[4]:

Text(0, 0.5, 'texture_mean')



In [5]:

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

Out[5]:

```
▼ KMeans  
KMeans()
```

In [6]:

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

Out[6]:

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

In [27]:

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

Out[27]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_me
0	842302	M	0.521037	0.022658	122.80	1001.0	0.118
1	842517	M	0.643144	0.272574	132.90	1326.0	0.084
2	84300903	M	0.601496	0.390260	130.00	1203.0	0.109
3	84348301	M	0.210090	0.360839	77.58	386.1	0.142
4	84358402	M	0.629893	0.156578	135.10	1297.0	0.100

5 rows × 35 columns

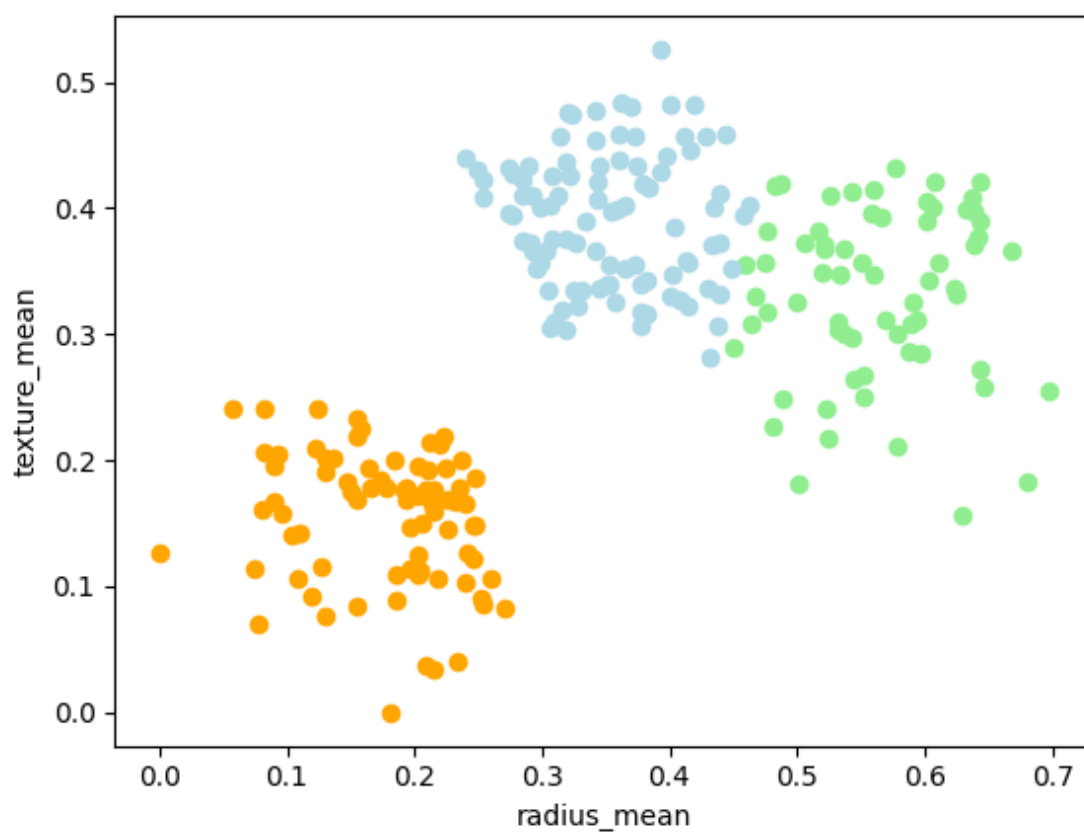


In [30]:

```
df1=df[df.Cluster==0]
df2=df[df.Cluster==1]
df3=df[df.Cluster==2]
plt.scatter(df1["radius_mean"],df1["texture_mean"],color="orange")
plt.scatter(df2["radius_mean"],df2["texture_mean"],color="lightgreen")
plt.scatter(df3["radius_mean"],df3["texture_mean"],color="lightblue")
plt.xlabel("radius_mean")
plt.ylabel("texture_mean")
```

Out[30]:

Text(0, 0.5, 'texture_mean')



In [9]:

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

Out[9]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_me
0	842302	M	17.99	0.022658	122.80	1001.0	0.118
1	842517	M	20.57	0.272574	132.90	1326.0	0.084
2	84300903	M	19.69	0.390260	130.00	1203.0	0.109
3	84348301	M	11.42	0.360839	77.58	386.1	0.142
4	84358402	M	20.29	0.156578	135.10	1297.0	0.100

5 rows × 33 columns

In [14]:

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

Out[14]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_me
0	842302	M	0.521037	0.022658	122.80	1001.0	0.118
1	842517	M	0.643144	0.272574	132.90	1326.0	0.084
2	84300903	M	0.601496	0.390260	130.00	1203.0	0.109
3	84348301	M	0.210090	0.360839	77.58	386.1	0.142
4	84358402	M	0.629893	0.156578	135.10	1297.0	0.100

5 rows × 33 columns

In [15]:

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

Out[15]:

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

In [16]:

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

Out[16]:

	id	diagnosis	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_me
0	842302	M	0.521037	0.022658	122.80	1001.0	0.118
1	842517	M	0.643144	0.272574	132.90	1326.0	0.084
2	84300903	M	0.601496	0.390260	130.00	1203.0	0.109
3	84348301	M	0.210090	0.360839	77.58	386.1	0.142
4	84358402	M	0.629893	0.156578	135.10	1297.0	0.100

5 rows × 34 columns

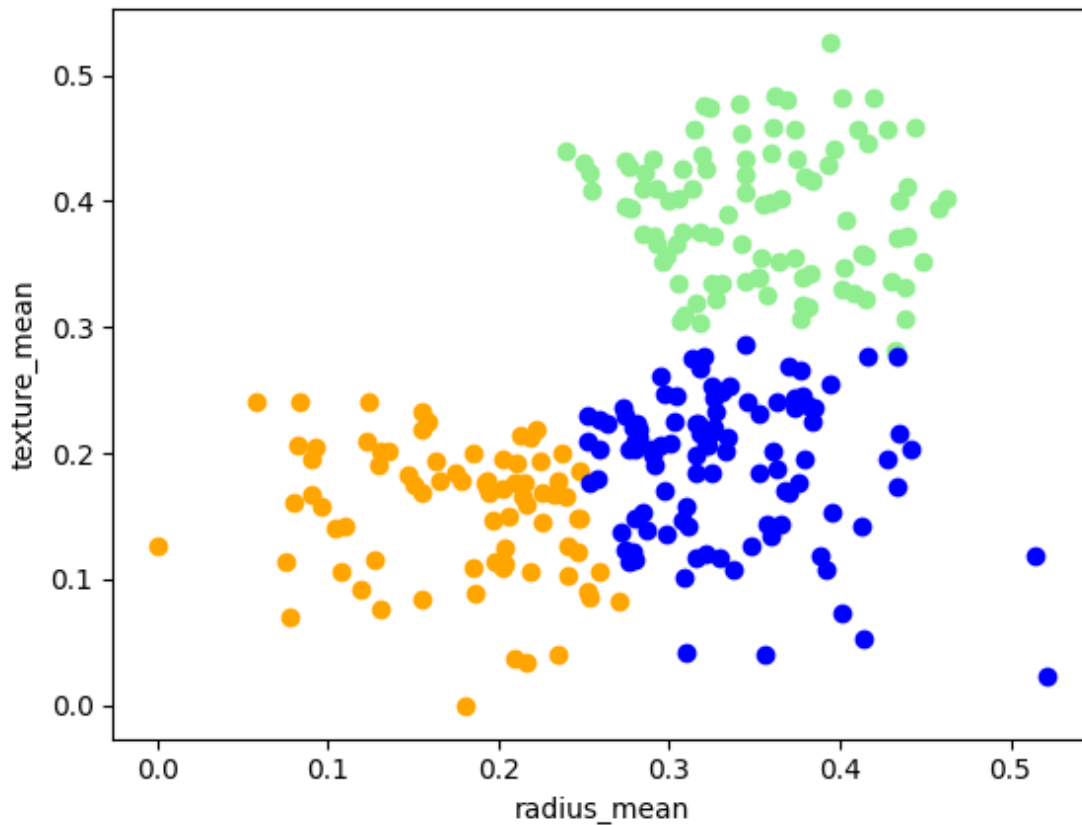


In [19]:

```
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="orange")
plt.scatter(df2["radius_mean"],df2["texture_mean"],color="lightgreen")
plt.scatter(df3["radius_mean"],df3["texture_mean"],color="blue")
plt.xlabel("radius_mean")
plt.ylabel("texture_mean")
```

Out[19]:

Text(0, 0.5, 'texture_mean')



In [20]:

```
km.cluster_centers_
```

Out[20]:

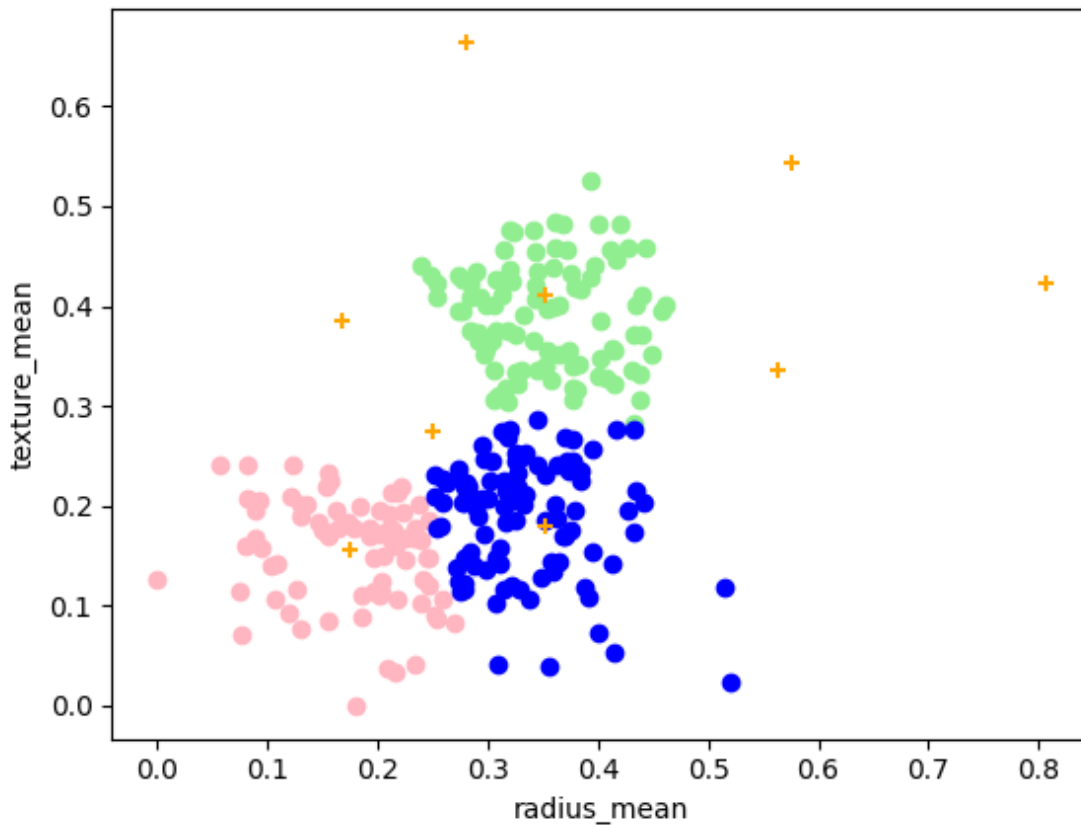
```
array([[0.17652977, 0.15382448],
       [0.56287997, 0.33184226],
       [0.3534653 , 0.39091896],
       [0.3331624 , 0.18999839],
       [0.57132058, 0.55893025],
       [0.20878924, 0.31058452],
       [0.2590623 , 0.58293879],
       [0.79840767, 0.42469846]])
```

In [31]:

```
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="lightpink")
plt.scatter(df2["radius_mean"],df2["texture_mean"],color="lightgreen")
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[31]:

Text(0, 0.5, 'texture_mean')



In [24]:

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

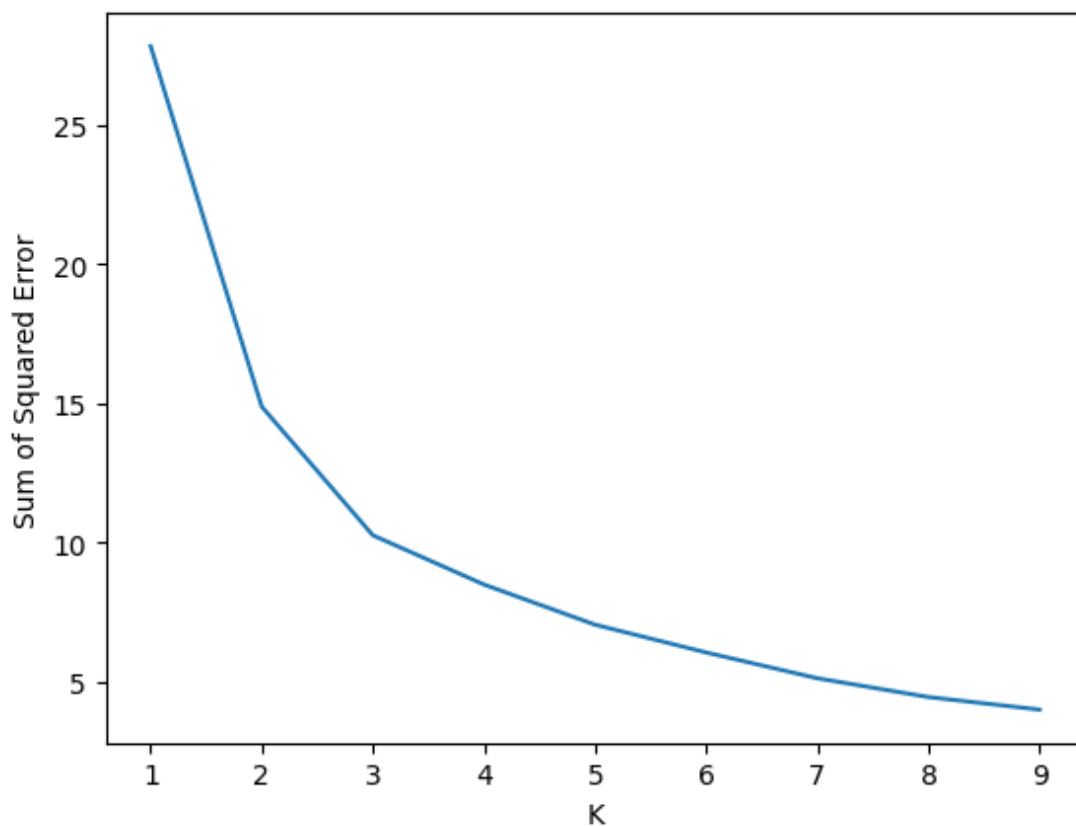

In [25]:

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

```
[27.817507595043075, 14.87203295827117, 10.252751496105198, 8.484720462531293,
7.04076556970335, 6.043008828607218, 5.117871158721308, 4.442425175619684, 3.9
921644950560777]
```

Out[25]:

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



Conclusion:-

In []:

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