

In []:

```
# Name : Meet Trivedi  
# Student Id; N01520331
```

Lab 9- K Mean Clustering

Import Libraries

In [1]:

```
import numpy as np  
import pandas as pd  
  
import matplotlib.pyplot as plt  
  
import sklearn  
from sklearn.preprocessing import scale  
import sklearn.metrics as sm  
from sklearn.metrics import confusion_matrix, classification_report  
from sklearn.cluster import KMeans  
from mpl_toolkits.mplot3d import Axes3D  
from sklearn import datasets  
%matplotlib inline
```

Get the Data(Wine dataset)

```
df = datasets.load_wine()
df
```

[illegible]

re thirteen different measurements taken for different constituents found in the three types of wine. Original Owners: Forina, M. et al, PARVUS - An Extendible Package for Data Exploration, Classification and Correlation. Institute of Pharmaceutical and Food Analysis and Technologies, Via Brigata Salerno, 16147 Genoa, Italy. Citation: Lichman, M. (2013). UCI Machine Learning Repository [https://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science. topic:: References

(1) S. Aeberhard, D. Coomans and O. de Vel, Comparison of Classifiers in High Dimensional Settings, Tech. Rep. no. 92-02, (1992), Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland. (Also submitted to Technometrics).

The data was used with many others for comparing various classifiers. The classes are separable, though only RDA has achieved 100% correct classification. (RDA : 100%, QDA 99.4%, LDA 98.9%, 1NN 96.1% (z-transformed data)) (All results using the leave-one-out technique)

(2) S. Aeberhard, D. Coomans and O. de Vel, "THE CLASSIFICATION PERFORMANCE OF RDA" Tech. Rep. no. 92-01, (1992), Dept. of Computer Science and Dept. of Mathematics and Statistics, James Cook University of North Queensland. (Also submitted to Journal of Chemometrics).

```
'feature_names': ['alcohol',
'malic_acid',
'ash',
'alcalinity_of_ash',
'magnesium',
'total_phenols',
'flavanoids',
'nonflavanoid_phenols',
'proanthocyanins',
'color_intensity',
'hue',
'od280/od315_of_diluted_wines',
'proline']}]
```

In [19]:

```
v_names = df.feature_names
v_names
```

Out[19]:

```
['alcohol',
'malic_acid',
'ash',
'alcalinity_of_ash',
'magnesium',
'total_phenols',
'flavanoids',
'nonflavanoid_phenols',
'proanthocyanins',
'color_intensity',
'hue',
'od280/od315_of_diluted_wines',
'proline']
```

In [20]:

```
x = scale(df.data)
x
```

Out[20]:

```
array([[ 1.51861254, -0.5622498 ,  0.23205254, ...,  0.36217728,
         1.84791957,  1.01300893],
       [ 0.24628963, -0.49941338, -0.82799632, ...,  0.40605066,
         1.1134493 ,  0.96524152],
       [ 0.19687903,  0.02123125,  1.10933436, ...,  0.31830389,
         0.78858745,  1.39514818],
       ...,
       [ 0.33275817,  1.74474449, -0.38935541, ..., -1.61212515,
        -1.48544548,  0.28057537],
       [ 0.20923168,  0.22769377,  0.01273209, ..., -1.56825176,
        -1.40069891,  0.29649784],
       [ 1.39508604,  1.58316512,  1.36520822, ..., -1.52437837,
        -1.42894777, -0.59516041]])
```

In [21]:

```
y = pd.DataFrame(df.target)
y
```

Out[21]:

	0
0	0
1	0
2	0
3	0
4	0
...	...
173	2
174	2
175	2
176	2
177	2

178 rows × 1 columns

Apply Kmean Clustering

In [34]:

```
clustering = KMeans(n_clusters=3, random_state=5)

clustering.fit(x)
```

Out[34]:

```
KMeans(n_clusters=3, random_state=5)
```

In [35]:

```
wine_df = pd.DataFrame(df.data)
wine_df.columns = ['alcohol', 'malic_acid', 'ash', 'alcalinity_of_ash', 'magnesium', 'total_
y.columns = ['Targets']
```

In [36]:

```
color_theme = np.array(['darkgray', 'lightsalmon', 'powderblue', 'darkblue', 'red', 'orange

plt.subplot(1,2,1)

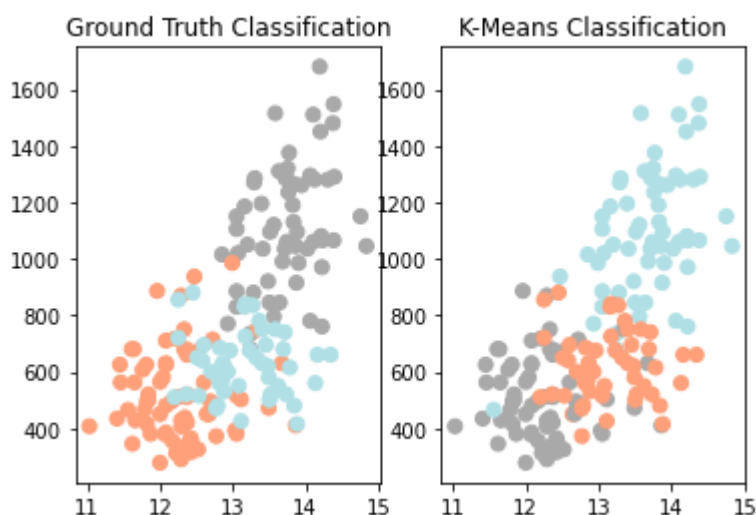
plt.scatter(x=wine_df.alcohol, y=wine_df.proline, c=color_theme[df.target], s=50)
plt.title('Ground Truth Classification')

plt.subplot(1,2,2)

plt.scatter(x=wine_df.alcohol, y=wine_df.proline, c=color_theme[clustering.labels_], s=50)
plt.title('K-Means Classification')
```

Out[36]:

```
Text(0.5, 1.0, 'K-Means Classification')
```



Predictions and Evaluations

In [38]:

```

relabel = np.choose(clustering.labels_, [2, 0, 1]).astype(np.int64)

plt.subplot(1,2,1)

plt.scatter(x=wine_df.alcohol, y=wine_df.proline, c=color_theme[df.target], s=50)
plt.title('Ground Truth Classification')

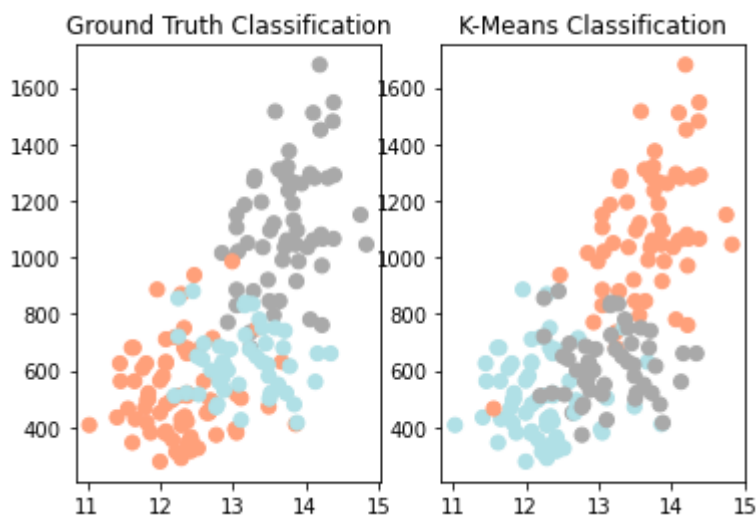
plt.subplot(1,2,2)

plt.scatter(x=wine_df.alcohol, y=wine_df.proline, c=color_theme[relabel], s=50)
plt.title('K-Means Classification')

```

Out[38]:

```
Text(0.5, 1.0, 'K-Means Classification')
```



In [39]:

```
print(classification_report(y, relabel))
```

	precision	recall	f1-score	support
0	0.00	0.00	0.00	59
1	0.05	0.04	0.05	71
2	0.00	0.00	0.00	48
accuracy			0.02	178
macro avg	0.02	0.01	0.02	178
weighted avg	0.02	0.02	0.02	178

In []:

