

## UnSupervised Learning Algorithm - Wine Color Prediction with K-Means

K-means is one of the most basic clustering algorithms. It relies on finding cluster centers to group data points based on minimizing the sum of squared errors between each datapoint and its cluster center. We are using a dataset which contains chemical properties (volatile\_acidity, total\_sulphur\_dioxide etc) to determine wine color

```
In [26]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import warnings
warnings.filterwarnings('ignore')
warnings.simplefilter('ignore')
```

```
In [27]: data = pd.read_csv("data/Wine_Quality_Data.csv")
```

```
In [28]: print(data.shape)
```

```
(6497, 13)
```

```
In [29]: #Print no of integers, floats and strings
data.dtypes.value_counts()
```

```
Out[29]: float64    11
int64         1
object        1
dtype: int64
```

```
In [30]: #Data should be numerical
data.head()
```

```
Out[30]:
```

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur
0	7.4	0.70	0.00	1.9	0.076	11.0	
1	7.8	0.88	0.00	2.6	0.098	25.0	
2	7.8	0.76	0.04	2.3	0.092	15.0	
3	11.2	0.28	0.56	1.9	0.075	17.0	
4	7.4	0.70	0.00	1.9	0.076	11.0	

```
In [31]: #Print no of entries for each color
data.color.value_counts()
```

```
Out[31]: white    4898
red        1599
Name: color, dtype: int64
```

```
In [32]: #Print % of each colors
data.color.value_counts(normalize=True)
```

```
Out[32]: white      0.753886
red          0.246114
Name: color, dtype: float64
```

## Preprocessing Steps

### 1. Select Features and apply feature tranformation/scaling.

```
In [33]: #Removing Color and Quality from features.
float_columns = [x for x in data.columns if x not in ['color', 'quality']]

# The correlation matrix
corr_mat = data[float_columns].corr()

#Every feature with itself will have correlation of one and we need to remove
for x in range(len(float_columns)):
    corr_mat.iloc[x,x] = 0.0

# max correlations(fixed_acidity,volatile_acidity has max co-relation )
corr_mat.abs().max()
```

```
Out[33]: fixed_acidity      0.458910
volatile_acidity      0.414476
citric_acid           0.377981
residual_sugar        0.552517
chlorides             0.395593
free_sulfur_dioxide   0.720934
total_sulfur_dioxide  0.720934
density              0.686745
pH                   0.329808
sulphates            0.395593
alcohol              0.686745
dtype: float64
```

```
In [34]: #Calculate Skew Vlaues
#0- no skew
#+ve - right skew
#-Ve - left skew
skew_columns = (data[float_columns]
                 .skew()
                 .sort_values(ascending=False))
skew_columns
```

```
Out[34]: chlorides          5.399828
sulphates          1.797270
fixed_acidity       1.723290
volatile_acidity    1.495097
residual_sugar      1.435404
free_sulfur_dioxide 1.220066
alcohol             0.565718
density             0.503602
citric_acid         0.471731
pH                  0.386839
total_sulfur_dioxide -0.001177
dtype: float64
```

```
In [35]: #Getting Skewed Columns and log tranforming it.
skew_columns = skew_columns.loc[skew_columns > 0.75]
# Perform log transform on skewed columns
for col in skew_columns.index.tolist():
    data[col] = np.log1p(data[col])
```

## 2. Normalize Features .

```
In [36]: from sklearn.preprocessing import StandardScaler
data[float_columns] = StandardScaler().fit_transform(data[float_columns])
data.head(4)
```

```
Out[36]:
```

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur
0	0.229509	2.135767	-2.192833	-0.815173	0.624554	-1.193601	-
1	0.550261	3.012817	-2.192833	-0.498175	1.281999	-0.013944	-
2	0.550261	2.438032	-1.917553	-0.625740	1.104012	-0.754684	-
3	2.802728	-0.337109	1.661085	-0.815173	0.594352	-0.574982	-

## Modeling with K-means

```
In [37]: from sklearn.cluster import KMeans
clusterNum = 2
k_means = KMeans(init = "k-means++", n_clusters = clusterNum, n_init = 12)
k_means.fit(data[float_columns])
labels = k_means.labels_
print(labels)

[1 1 1 ... 0 0 0]
```

```
In [38]: #Assigning labels generated by K-means to our original dataset
data["Kmeans"] = labels
data.head(5)
```

Out[38]:

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur
0	0.229509	2.135767	-2.192833	-0.815173	0.624554	-1.193601	-
1	0.550261	3.012817	-2.192833	-0.498175	1.281999	-0.013944	-
2	0.550261	2.438032	-1.917553	-0.625740	1.104012	-0.754684	-
3	2.802728	-0.337109	1.661085	-0.815173	0.594352	-0.574982	-
4	0.229509	2.135767	-2.192833	-0.815173	0.624554	-1.193601	-

```
In [39]: data.groupby('Kmeans').mean()
```

Out[39]:

	fixed_acidity	volatile_acidity	citric_acid	residual_sugar	chlorides	free_sulfur_dioxide	total_sulfur
Kmeans							
0	-0.276629	-0.394800	0.111070	0.202578	-0.328822	0.336927	
1	0.804104	1.147603	-0.322858	-0.588853	0.955818	-0.979379	

```
In [40]: #Without giving labels K-Means has created two clusters and lets examine how
(data[['color', 'Kmeans']]
 .groupby(['Kmeans', 'color'])
 .size()
 .to_frame()
 .rename(columns={0: 'number'}))
```

Out[40]:

		number
Kmeans	color	
0	red	23
	white	4811
1	red	1576
	white	87

**Kmeans could classify into two clusters one with having red as majority and another with**

**having white has majority.**

## How to use a inertia Curve to determine optimal number of Clusters?

In [41]:

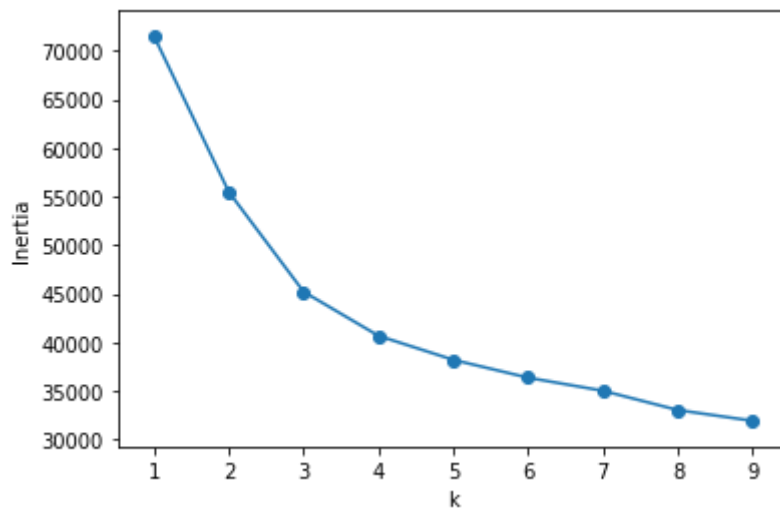
```
k_range = range(1,10)
inertia = []
for k in k_range:
    k_means = KMeans(init = "k-means++", n_clusters = k, n_init = 12)
    k_means.fit(data[float_columns])
    inertia.append(k_means.inertia_)
```

In [42]: `print(inertia)`

```
[71467.0, 55455.9643286125, 45199.872248980435, 40674.02951127552, 38203.
86301123881, 36367.79026347645, 35023.58208190495, 33040.88734181203, 319
34.672934147253]
```

In [43]:

```
plt.plot(k_range,inertia)
plt.scatter(k_range,inertia)
plt.xlabel('k')
plt.ylabel('Inertia');
```



## Summary

***Without providing labels K-means can classify wine to red and white with one cluster having majority red and another with white .***

Elbow Curve helps us to determine the number of Clusters required . Inertia continues to go down as number of clusters increases but after sometime number it flattens down. Here we can choose number of clusters as 4. Since we know that there are only two values it is better to choose two clusters.

