

## Import Necessary Libraries

```
In [39]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.decomposition import PCA
from sklearn.preprocessing import scale
```

## Business Problem

Performing Principal component analysis and clustering using first 3 principal component scores (both heirarchical and k mean clustering) and obtain optimum number of clusters and check whether we have obtained same number of clusters with the original data (class column we have ignored at the begining who shows it has 3 clusters)df

## Data Understanding

In [31]:

wine=pd.read_csv('wine.csv')													
wine													
Type	Alcohol	Malic	Ash	Alcalinity	Magnesium	Phenols	Flavanoids	Nonflavanoids	Proanthocyanins	Color	Hue	Dilution	Proline
0	1	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92
1	1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40
2	1	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17
3	1	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45
4	1	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93
...	...	...	...	...	...	...	...	...	...	...	...	...	...
173	3	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.70	0.64	1.74
174	3	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.30	0.70	1.56
175	3	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.20	0.59	1.56
176	3	13.17	2.59	2.37	20.0	120	1.65	0.68	0.53	1.46	9.30	0.60	1.62
177	3	14.13	4.10	2.74	24.5	96	2.05	0.76	0.56	1.35	9.20	0.61	1.60

178 rows × 14 columns

## Data Preparation

Initial Analysis

In [32]:

wine['Type'].value_counts()													
2	71												
1	59												
3	48												
Name: Type, dtype: int64													

Out [32]:

In [33]:

wine_data=wine.iloc[:,1:]													
wine_data													
Alcohol	Malic	Ash	Alcalinity	Magnesium	Phenols	Flavanoids	Nonflavanoids	Proanthocyanins	Color	Hue	Dilution	Proline	
0	14.23	1.71	2.43	15.6	127	2.80	3.06	0.28	2.29	5.64	1.04	3.92	1065
1	13.20	1.78	2.14	11.2	100	2.65	2.76	0.26	1.28	4.38	1.05	3.40	1050
2	13.16	2.36	2.67	18.6	101	2.80	3.24	0.30	2.81	5.68	1.03	3.17	1185
3	14.37	1.95	2.50	16.8	113	3.85	3.49	0.24	2.18	7.80	0.86	3.45	1480
4	13.24	2.59	2.87	21.0	118	2.80	2.69	0.39	1.82	4.32	1.04	2.93	735
...	...	...	...	...	...	...	...	...	...	...	...	...	...
173	13.71	5.65	2.45	20.5	95	1.68	0.61	0.52	1.06	7.70	0.64	1.74	740
174	13.40	3.91	2.48	23.0	102	1.80	0.75	0.43	1.41	7.30	0.70	1.56	750
175	13.27	4.28	2.26	20.0	120	1.59	0.69	0.43	1.35	10.20	0.59	1.56	835
176	13.17	2.59	2.37	20.0	120	1.65	0.68	0.53	1.46	9.30	0.60	1.62	840
177	14.13	4.10	2.74	24.5	96	2.05	0.76	0.56	1.35	9.20	0.61	1.60	560

178 rows × 13 columns

Out [33]:

In [34]:

wine_data.shape													
(178, 13)													

Out [34]:

In [35]:

wine_data.info()													
<class 'pandas.core.frame.DataFrame'>													
RangeIndex: 178 entries, 0 to 177													
Data columns (total 13 columns):													
# Column Non-Null Count Dtype													
---													
0	Alcohol	178 non-null	float64										
1	Malic	178 non-null	float64										
2	Ash	178 non-null	float64										
3	Alcalinity	178 non-null	float64										
4	Magnesium	178 non-null	int64										
5	Phenols	178 non-null	float64										
6	Flavanoids	178 non-null	float64										
7	Nonflavanoids	178 non-null	float64										
8	Proanthocyanins	178 non-null	float64										
9	Color	178 non-null	float64										
10	Hue	178 non-null	float64										
11	Dilution	178 non-null	float64										
12	Proline	178 non-null	int64										
dtypes: float64(11), int64(2)													
memory usage: 18.2 KB													

Out [35]:

In [36]:

# Converting data to numpy array													
wine_ary=wine_data.values													
wine_ary													
array([[ 1.423e+01,  1.710e+00,  2.430e+00, ...,  1.040e+00,  3.920e+00,													
[ 1.065e+03],													
[ 1.320e+01,  1.780e+00,  2.140e+00, ...,  1.050e+00,  3.400e+00,													
[ 1.050e+03],													
[ 1.316e+01,  2.360e+00,  2.670e+00, ...,  1.030e+00,  3.170e+00,													
[ 1.185e+03],													
...,													
[ 1.327e+01,  4.280e+00,  2.260e+00, ...,  5.900e-01,  1.560e+00,													
[ 8.350e+02],													
[ 1.317e+01,  2.590e+00,  2.370e+00, ...,  6.000e-01,  1.620e+00,													
[ 8.400e+02],													
[ 1.413e+01,  4.100e+00,  2.740e+00, ...,  6.100e-01,  1.600e+00,													
[ 5.600e+02]])													

Out [36]:

In [37]:

# Normalizing the numerical data													
wine_norm=scale(wine_ary)													
wine_norm													
array([[ 1.51861254, -0.5622498,  0.23205254, ...,  0.36217728,													
[ 1.84791957,  1.01308093],													
[ 0.24628963, -0.4994138, -0.82799632, ...,  0.40605966,													
[ 1.134403,  0.06504152],													
[ 0.19687903,  0.02123125,  1.10933436, ...,  0.31830389,													
[ 0.78858745,  1.39514818],													
...													
[ 0.33275817,  1.74474449, -0.38935541, ..., -1.61212515,													
[ -1.40544548,  0.20057537],													
[ 0.20923108,  0.27708377,  0.01273209, ..., -1.56825176,													
[ -1.40069891,  0.29649784],													
[ 1.39508094,  1.58316512,  1.36520622, ..., -1.52437837,													
[ -1.42047777,  0.59516041]])													

Out [37]:

## PCA Implementation

```
# -2.38761709e+00, -2.29734668e+00, ...,
2.99821968e-01, 3.39820654e-01, -2.18657665e-02],
[-3.29875816e+00, 2.76801907e+00, 1.01301366e+00, ...,
-2.29964331e-01, -1.88787963e-01, -3.23964720e-01]])

# PCA Components matrix or covariance Matrix
pca.components_

array([[ 0.1443294, -0.24518758, -0.00205106, -0.23932041, 0.14199204,
0.39466085, 0.4229343, -0.2985331, 0.31342049, -0.0886167,
0.296741556, 0.37618741, 0.20675222],
[ -0.48365155, -0.22493893, 0.31606881, 0.0105095, -0.299634,
-0.06503951, 0.00335981, 0.02877949, -0.03930172, -0.52999567,
0.27923515, 0.16440813, 0.36408263],
[ -0.20728262, 0.08901289, 0.6262239, 0.61208035, 0.13075693,
0.14617896, 0.1560819, 0.17036816, 0.14945431, -0.13730621,
0.08521392, 0.16680459, 0.12674592],
[ -0.0178563, 0.53689828, -0.21417556, 0.06085941, -0.35179658,
0.19096835, 0.15229479, -0.20330102, 0.39905653, -0.13769621,
0.427717141, 0.18412074, -0.23207698],
[ -0.26563655, 0.03521363, -0.14302547, 0.06610294, 0.72704851,
-0.14931841, -0.10902584, -0.50070298, 0.13685982, -0.07643078,
0.17361452, -0.10116999, 0.15716868],
[ -0.21353865, -0.53681385, -0.15447466, 0.10082451, -0.03814394,
0.0641223, 0.01892892, -0.25859461, 0.53379539, 0.41864414,
-0.10598274, -0.26585107, -0.11972557],
[ -0.05639636, 0.42652391, -0.14917061, -0.28696914, 0.3228833,
-0.02724498, -0.06068521, 0.59547729, 0.37213935, -0.22771214,
0.23207564, -0.0447637, 0.0769845],
[ -0.39613926, -0.06582674, 0.17026002, -0.42797018, 0.15636143,
0.40593409, 0.18724536, 0.23328465, -0.36822075, 0.83379692,
-0.43662362, 0.07810789, -0.12062267],
[ 0.50861912, -0.07528304, 0.30769445, 0.20844931, 0.27140257,
0.28603452, 0.04957849, 0.19560132, -0.20914487, 0.06621752,
0.08582839, 0.1372269, -0.57578611],
[ 0.21160473, -0.30907994, -0.02712539, 0.05279942, 0.06787022,
0.32013135, -0.16310951, -0.21553507, 0.1341839, -0.28075218,
-0.52239889, 0.52370587, 0.162116],
[ -0.22591696, 0.07648554, -0.49869142, 0.47931378, 0.07128991,
0.38434119, -0.02569489, 0.11689596, -0.23736257, 0.031388,
-0.04821201, 0.0464233, -0.53926983],
[ -0.26286455, 0.12169684, -0.04962237, -0.05674287, 0.06222011,
-0.30388245, -0.04289883, 0.04235219, -0.09555393, 0.06421163,
0.259214, 0.606995872, -0.07940162],
[ 0.01409097, 0.02595375, -0.14121803, 0.09168285, 0.05677422,
-0.46309791, 0.83225706, 0.11403985, -0.11691707, -0.0119928,
-0.08988884, -0.15671813, 0.01444734]])

# The amount of variance that each PCA has
var=np.pca.explained_variance_ratio_
var

array([0.36198848, 0.1920749, 0.11123631, 0.0706903, 0.06563294,
0.04935823, 0.04238679, 0.02680749, 0.0222153, 0.01930019,
0.01736836, 0.01298233, 0.00795215])

# Cumulative variance of each PCA
var=np.cumsum(np.round(var,4)*100)
var1

array([ 36.2, 55.41, 66.53, 73.6, 80.16, 85.1, 89.34, 92.82,
94.24, 96.17, 97.01, 99.21, 100.01])

# Variance plot for PCA components obtained
plt.plot(var1,color='magenta')

[<matplotlib.lines.Line2D at 0x287b7418dc0>]

100
80
60
40
0
0 2 4 6 8 10 12

# Final DataFrame
final_df=pd.concat([wine['Type'],pd.DataFrame(wine_pca[:,0:3],columns=['PC1','PC2','PC3'])],axis=1)
final_df

Type PC1 PC2 PC3
0 1 3.316751 -1.443403 -0.165739
1 1 2.209465 0.333393 -0.026457
2 1 2.516740 -1.031151 0.982819
3 1 3.757066 -2.756372 -0.176192
4 1 1.008908 -0.869831 2.026688
... ..
173 3 -3.370524 -2.216289 -0.342570
174 3 -2.601956 -1.757229 0.207581
175 3 -2.677839 -2.760899 -0.940942
176 3 -2.387017 -2.297347 -0.550696
177 3 -3.208758 -2.768920 1.013914

178 rows x 4 columns

# Visualization of PCAs
fig=plt.figure(figsize=(16,12))
sns.scatterplot(data=final_df)
```