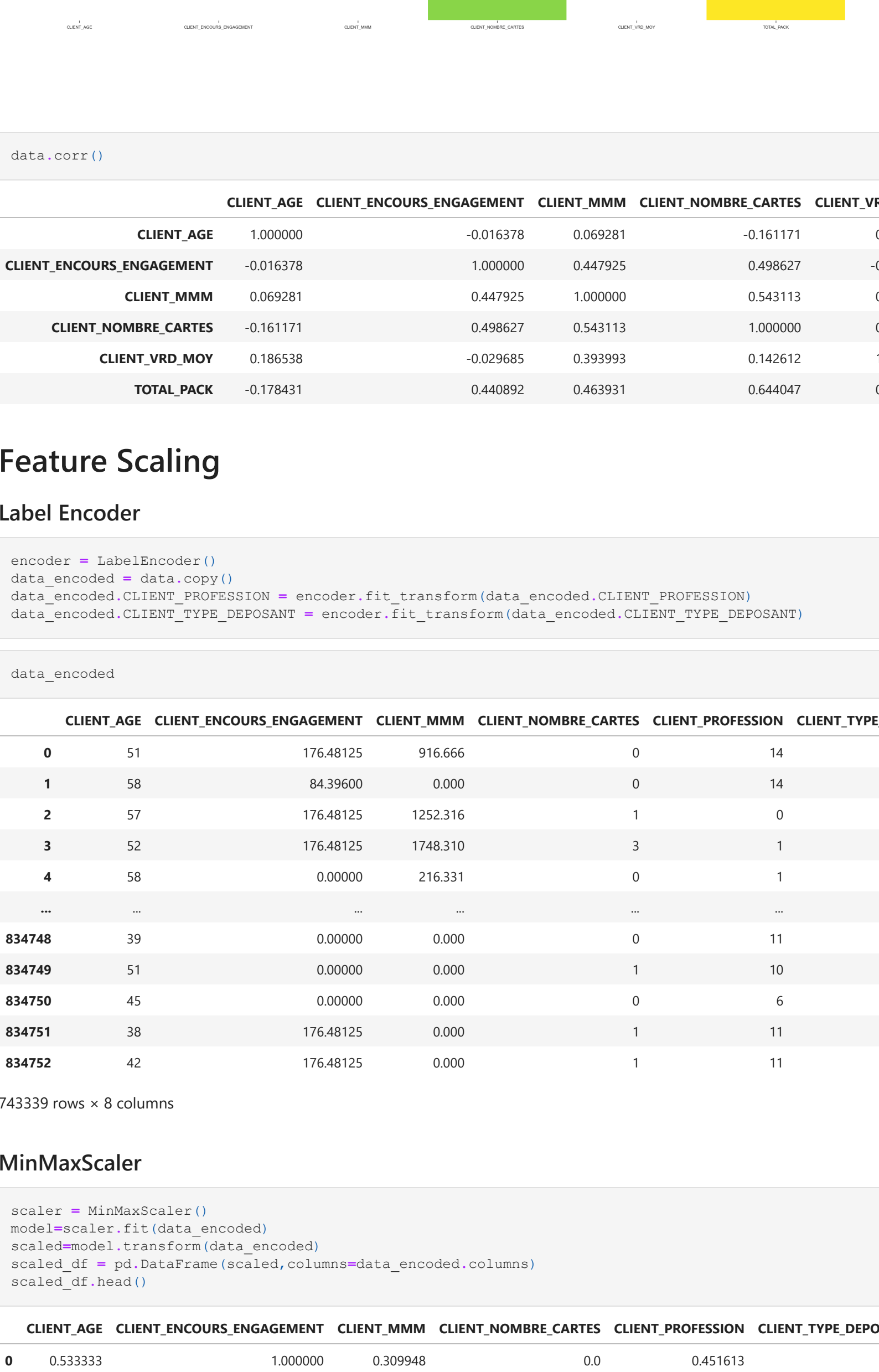


```
plt.figure(figsize=(50, 50))
sns.heatmap(correlation[correlation >= 0.5 | (correlation <= -0.4)],
            cmap='viridis', vmax=1.0, vmin=-1.0, linewidths=0.1,
            annot=True, annot_kws={"size": 8}, square=True);
```



	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	TOT
CLIENT_AGE	1.000000	-0.016378	0.069281	-0.161171			0.186538
CLIENT_ENCOURS_ENGAGEMENT	-0.016378	1.000000	0.447925	0.498627			-0.029685
CLIENT_MMM	0.069281	0.447925	1.000000	0.543113			0.393993
CLIENT_NOMBRE_CARTES	-0.161171	0.498627	0.543113	1.000000			0.142612
CLIENT_VRD_MOY	0.186538	-0.029685	0.393993	0.142612			1.000000
TOTAL_PACK	-0.178431	0.440892	0.463931	0.644047			0.132922

## Feature Scaling

### Label Encoder

```
encoder = LabelEncoder()
data_encoded = data.copy()
data_encoded.CLIENT_PROFESSION = encoder.fit_transform(data_encoded.CLIENT_PROFESSION)
data_encoded.CLIENT_TYPE_DEPOSANT = encoder.fit_transform(data_encoded.CLIENT_TYPE_DEPOSANT)
```

```
data_encoded.head()
```

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	TOT
0	51	176.48125	916.666	0	14		0
1	58	84.39600	0.000	0	14		0
2	57	176.48125	1252.316	1	0		0
3	58	176.48125	1748.310	3	1		0
4	52	0.00000	216.331	0	1		0
...	...	...	...	...	...	...	...
834748	39	0.00000	0.000	0	11		1
834749	51	0.00000	0.000	1	10		4
834750	48	0.00000	0.000	0	6		1
834751	35	176.48125	0.000	1	11		0
834752	42	176.48125	0.000	1	11		0

743339 rows x 8 columns

### MinMaxScaler

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	TOT
0	0.333333	1.000000	0.309948	0.0	0.451613		0.0
1	0.650000	0.478215	0.000000	0.0	0.451613		0.0
2	0.633333	1.000000	0.423440	0.2	0.000000		0.0
3	0.550000	1.000000	0.591148	0.6	0.032258		0.0
4	0.650000	0.000000	0.073147	0.0	0.032258		0.0

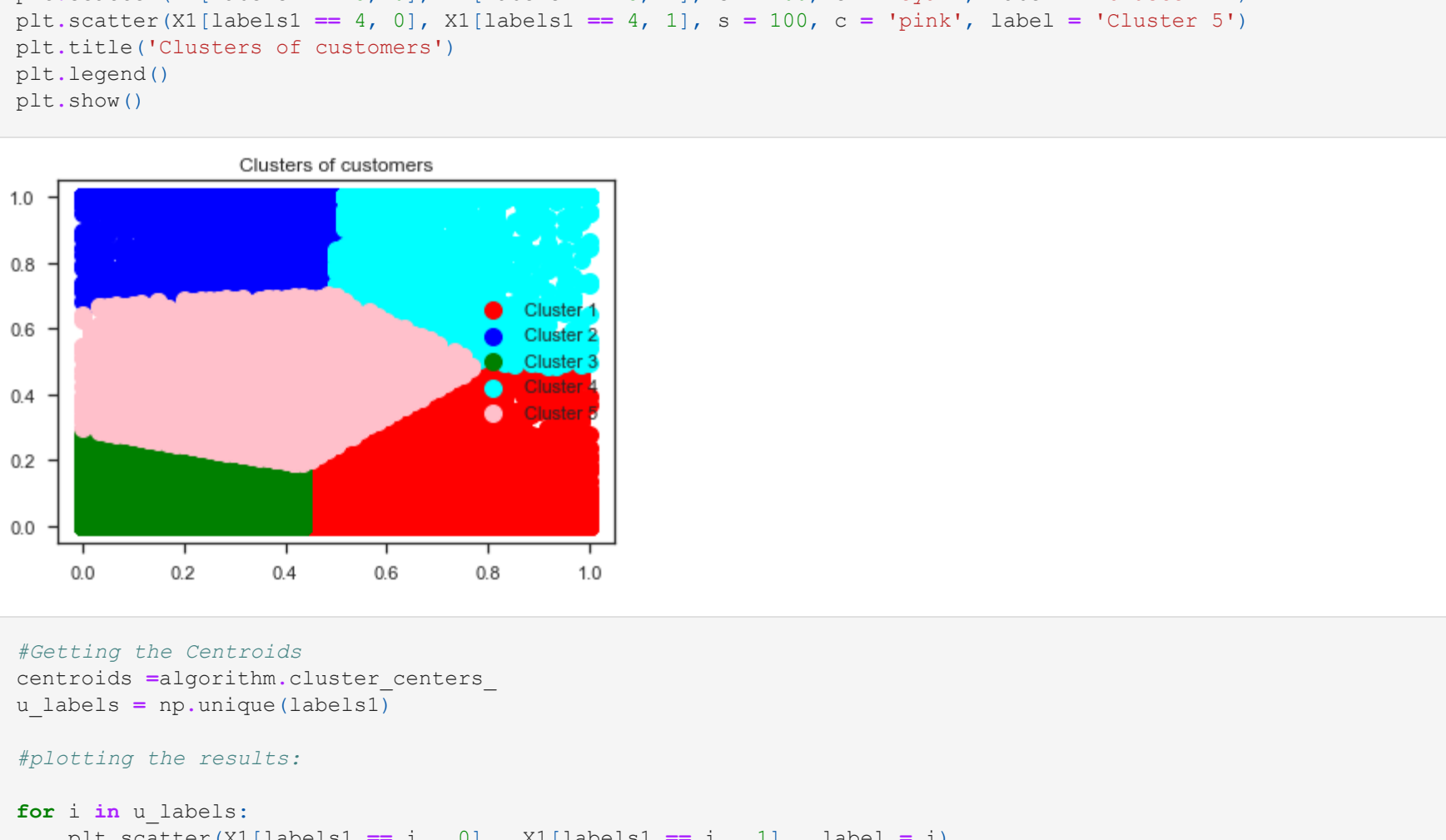
## Cluster Analysis

### Elbow method

```
temp = scaled_df.sample(n=50000)
```

```
X1 = temp[['CLIENT_AGE', 'CLIENT_VRD_MOY']].iloc[:, :].values
inertia = []
s_scores = []
for i in range(2, 11):
    algorithm = KMeans(n_clusters = n, init='k-means++', n_init = 10, max_iter=300,
                        tol=0.0001, random_state=111, algorithm='elkan')
    algorithm.fit(X1)
    inertia.append(algorithm.inertia_)
    silhouette_avg = silhouette_score(X1, algorithm.labels_)
    s_scores.append(silhouette_avg) # data for the silhouette score method
```

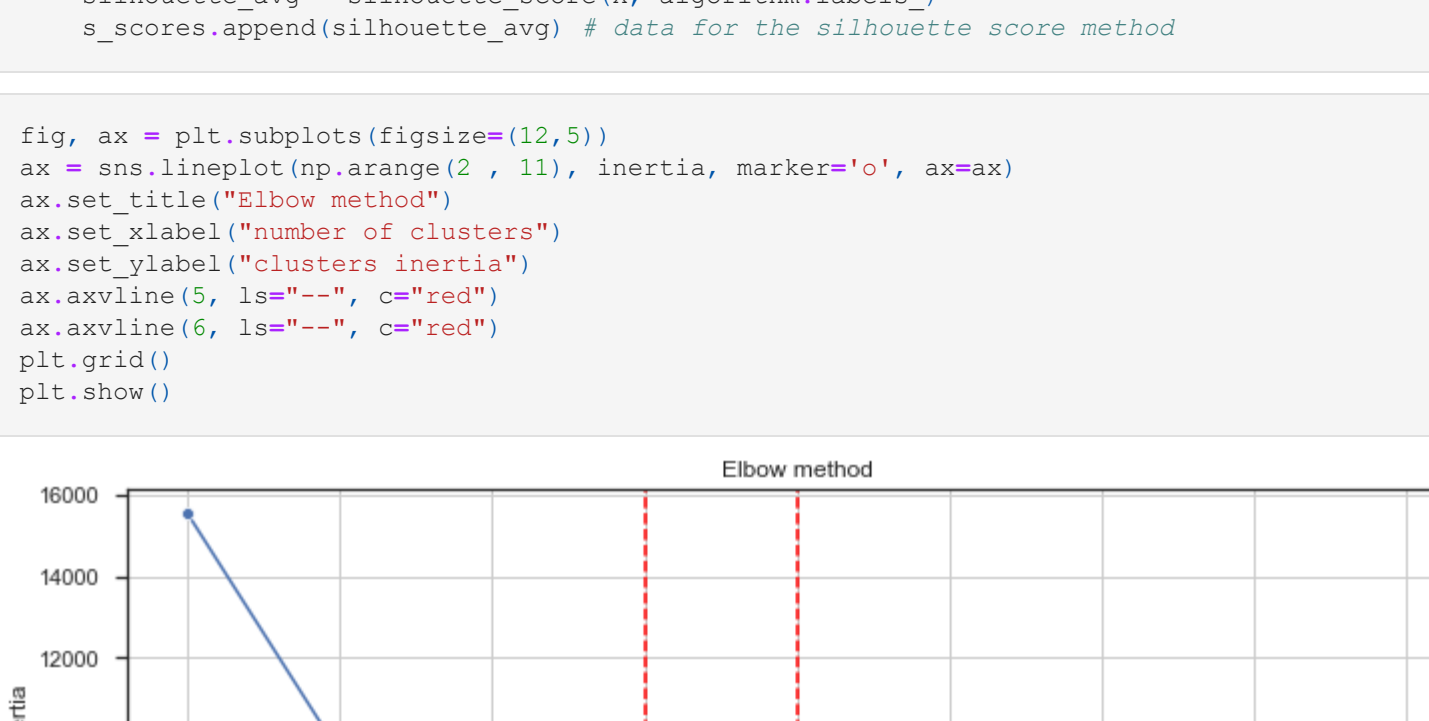
```
#The Elbow Chart
plt.figure(1, figsize=(15, 6))
plt.plot(np.arange(2, 11), inertia, 'o')
plt.plot(np.arange(2, 11), inertia, '-', alpha=0.5)
plt.xlabel('Number of Clusters',)
plt.ylabel('Inertia')
```



It is unclear which clusters to go for 4 or 5 or 6. Here we use the Silhouette method. The range of the Silhouette value is between +1 and -1. A high value is desirable and indicates that the point is placed in the correct cluster. If many points have a negative Silhouette value, it may indicate that we have created too many or too few clusters.

A high Silhouette Score is desirable. The Silhouette Score reaches its global maximum at the optimal k. This should ideally appear as a peak in the Silhouette Value-versus-k plot.

```
fig, ax = plt.subplots(figsize=(12,5))
ax = sns.lineplot(np.arange(2, 11), s_scores, marker='o', ax=ax)
ax.set_title('Silhouette score method')
ax.set_xlabel('Number of clusters')
ax.set_ylabel('Silhouette score')
ax.axvline(5, ls='--', c='red')
ax.axvline(6, ls='--', c='red')
ax.grid()
plt.show()
```

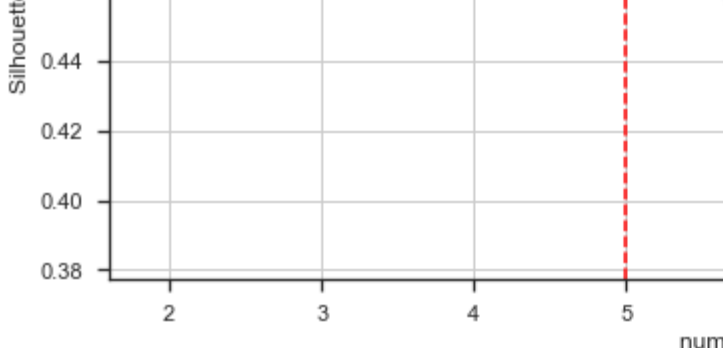


We see that silhouette score is maximum at k=5. So we go ahead and take that value

The Elbow Method is more of a decision rule, while the Silhouette is a metric used for validation while clustering. Thus, it can be used in combination with the Elbow Method.

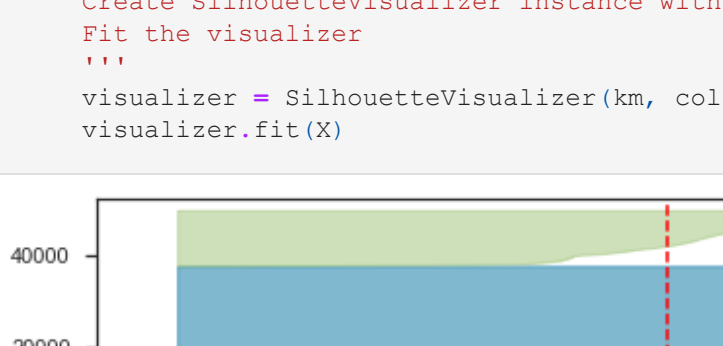
```
algorithm = KMeans(n_clusters = 5, init='k-means++', n_init = 10, max_iter=300,
                    tol=0.0001, random_state=111, algorithm='elkan')
algorithm.fit(X1)
labels1 = algorithm.labels_
cluster_centers1 = algorithm.cluster_centers_
```

```
#Plotting the clusters
plt.scatter(X1[labels1 == 0, 0], X1[labels1 == 0, 1], s = 100, c = 'red', label = 'Cluster 1')
plt.scatter(X1[labels1 == 1, 0], X1[labels1 == 1, 1], s = 100, c = 'blue', label = 'Cluster 2')
plt.scatter(X1[labels1 == 2, 0], X1[labels1 == 2, 1], s = 100, c = 'green', label = 'Cluster 3')
plt.scatter(X1[labels1 == 3, 0], X1[labels1 == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4')
plt.scatter(X1[labels1 == 4, 0], X1[labels1 == 4, 1], s = 100, c = 'pink', label = 'Cluster 5')
plt.legend()
plt.show()
```



```
#Getting the Centroids
centroids = algorithm.cluster_centers_
u_labels = np.unique(labels1)

#Plotting the results:
for i in u_labels:
    plt.scatter(X1[labels1 == i, 0], X1[labels1 == i, 1], label = i)
plt.scatter(centroids[i, 0], centroids[i, 1], s = 80, color = 'k')
plt.legend()
plt.show()
```



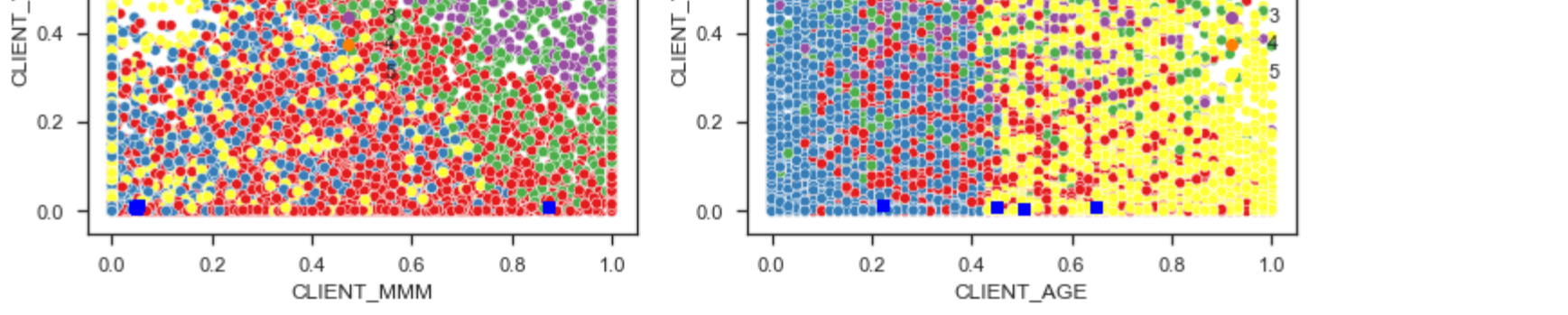
### Performing segmentation using all the numeric variables

```
X = temp[['CLIENT_AGE', 'CLIENT_MMM', 'CLIENT_ENCOURS_ENGAGEMENT', 'CLIENT_VRD_MOY', 'TOTAL_PACK']]
inertia = []
s_scores = []
for i in range(2, 11):
    algorithm = KMeans(n_clusters = n, init='k-means++', n_init = 10, max_iter=300,
                        tol=0.0001, random_state=111, algorithm='elkan').fit(X)
    inertia.append(algorithm.inertia_)
    silhouette_avg = silhouette_score(X, algorithm.labels_)
    s_scores.append(silhouette_avg) # data for the silhouette score method
```

```
fig, ax = plt.subplots(figsize=(12,5))
ax = sns.lineplot(np.arange(2, 11), s_scores, marker='o', ax=ax)
ax.set_title('Elbow method')
ax.set_xlabel('Number of clusters')
ax.set_ylabel('Silhouette score')
ax.axvline(5, ls='--', c='red')
ax.axvline(6, ls='--', c='red')
plt.grid()
plt.show()
```



```
fig, ax = plt.subplots(3, 2, figsize=(15,8))
for i in [2,3,4,5,6,7]:
    km = KMeans(n_clusters=i, init='k-means++', n_init=10, max_iter=300, random_state=42)
    km.fit(X)
    mod = km.predict(X)
    visualizer = SilhouetteVisualizer(km, colors='yellowbrick', ax=ax[i-1][mod])
    visualizer.fit(X)
```



```
#For clusters of K=6
KM6 = KMeans(n_clusters = 6, init='k-means++', n_init = 10, max_iter=300,
              tol=0.0001, random_state=111, algorithm='elkan')
KM6.fit(X)
labels6 = KM6.labels_
centroids6 = KM6.cluster_centers_
KM6_df = temp.copy()
KM6_df['labels'] = labels6
```

```
fig1, (axes) = plt.subplots(1,2,figsize=(12,5))

scat1 = sns.scatterplot('CLIENT_MMM', 'CLIENT_VRD_MOY', data=KM6_df,
                        hue='labels', ax=axes[0], palette='Set1', legend='full')

sns.scatterplot('CLIENT_AGE', 'CLIENT_VRD_MOY', data=KM6_df,
                hue='labels', palette='Set1', ax=axes[1], legend='full')
```

```
axes[0].scatter(centroids6[1,1],centroids6[1,2], marker='s', s=40, c='blue')
axes[1].scatter(centroids6[1,0],centroids6[1,2], marker='s', s=40, c='blue')
plt.show()
```



```
KM_clust_sizes = KM6_df.groupby('labels').size().to_frame()
KM_clust_sizes.columns = ['KM_size']
KM_clust_sizes
```

	KM_size
0	9469
1	17963
2	4336
3	2765
4	5562
5	9905

```
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure(figsize=(7,7))
ax = Axes3D(fig, rect=(0, 0, .99, 1), elev=20, azim=210)
ax.scatter(KM6_df['CLIENT_AGE'],
           KM6_df['CLIENT_MMM'],
           KM6_df['CLIENT_VRD_MOY'],
           c=KM6_df['labels'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)
```

```
ax.w_axis.set_ticklabels([])
ax.y_axis.set_ticklabels([])
ax.z_axis.set_ticklabels([])
ax.set_xlabel('CLIENT_AGE')
ax.set_ylabel('CLIENT_MMM')
ax.set_title('3D view of K-Means 6 clusters')
ax.dist = 12
plt.show()
```



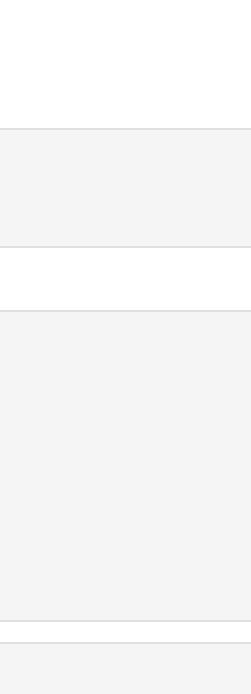
```
import plotly as py
import plotly.graph_objs as go
def tracer(db, n, name):
    This function returns trace object for Plotly
    return go.Scatter3d(x=db[db['labels']==n]['CLIENT_AGE'],
                        y=db[db['labels']==n]['CLIENT_VRD_MOY'],
                        z=db[db['labels']==n]['CLIENT_MMM'],
                        mode = 'markers',
                        name = name,
                        marker = dict(
                            color = KM6_df['labels'],
                            size = 5,
                            linedict = KM6_df['labels'],
                            width= 12
                        ),
                        opacity=0.8
                    )
```

```
trace0 = tracer(KM6_df, 0, 'Cluster 0')
trace1 = tracer(KM6_df, 1, 'Cluster 1')
trace2 = tracer(KM6_df, 2, 'Cluster 2')
trace3 = tracer(KM6_df, 3, 'Cluster 3')
trace4 = tracer(KM6_df, 4, 'Cluster 4')
trace5 = tracer(KM6_df, 5, 'Cluster 5')
trace_data = [trace0, trace1, trace2, trace3, trace4, trace5]
```

```
layout = go.Layout(
    title = 'Clusters wrt Age, Income and Spending Scores',
    scene = dict(
        xaxis = dict(title = 'Age'),
        yaxis = dict(title = 'Spending Score'),
        zaxis = dict(title = 'Annual Income')
    )
)
```

```
fig = go.Figure(data=trace_data, layout=layout)
py.offline.iplot(fig)
```

Clusters with k=6 wrt Age, Income and Spending Scores



```
db_index = davis_bouldin_score(X, labels6)
print('Daves Bouldin Index: %.2f'%db_index)
```

Daves Bouldin Index: 0.92

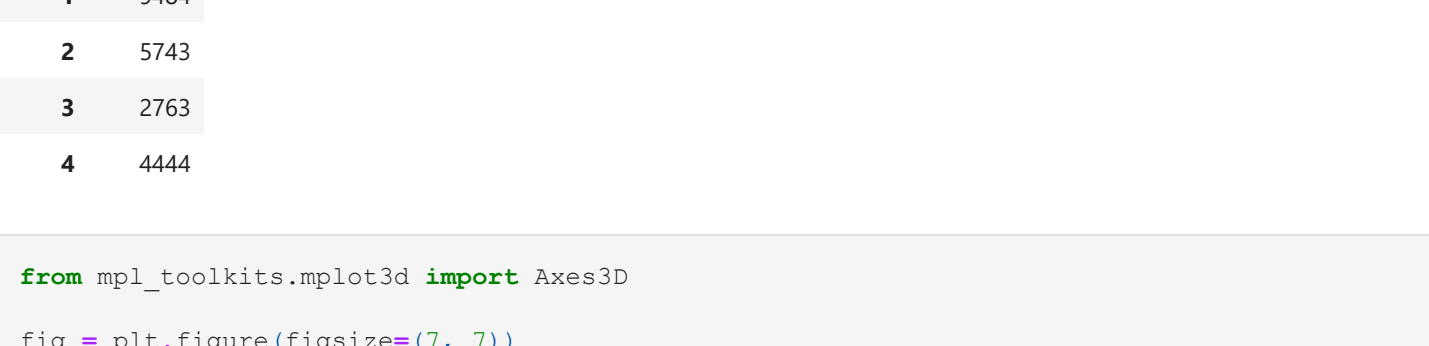
```
#For clusters of K=5
KM5 = KMeans(n_clusters = 5, init='k-means++', n_init = 10, max_iter=300,
              tol=0.0001, random_state=111, algorithm='elkan')
KM5.fit(X)
labels5 = KM5.labels_
centroids5 = KM5.cluster_centers_
KM5_df = temp.copy()
KM5_df['labels'] = labels5
```

```
fig1, (axes) = plt.subplots(1,2,figsize=(12,5))

scat1 = sns.scatterplot('CLIENT_MMM', 'CLIENT_VRD_MOY', data=KM5_df,
                        hue='labels', ax=axes[0], palette='Set1', legend='full')

sns.scatterplot('CLIENT_AGE', 'CLIENT_VRD_MOY', data=KM5_df,
                hue='labels', palette='Set1', ax=axes[1], legend='full')
```

```
axes[0].scatter(centroids5[1,1],centroids5[1,2], marker='s', s=40, c='blue')
axes[1].scatter(centroids5[1,0],centroids5[1,2], marker='s', s=40, c='blue')
plt.show()
```

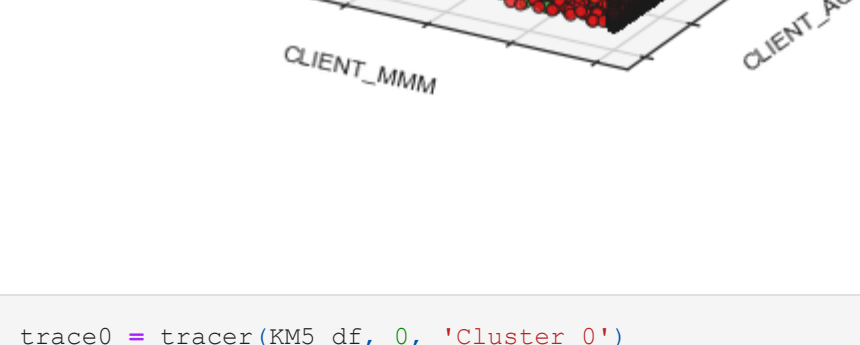


```
KM_clust_sizes5 = KM5_df.groupby('labels').size().to_frame()
KM_clust_sizes5.columns = ['KM_size']
KM_clust_sizes5
```

	KM_size
0	27566
1	9484
2	5743
3	2763
4	4444

```
from mpl_toolkits.mplot3d import Axes3D
fig = plt.figure(figsize=(7,7))
ax = Axes3D(fig, rect=(0, 0, .99, 1), elev=20, azim=210)
ax.scatter(KM5_df['CLIENT_AGE'],
           KM5_df['CLIENT_MMM'],
           KM5_df['CLIENT_VRD_MOY'],
           c=KM5_df['labels'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)
```

```
ax.w_axis.set_ticklabels([])
ax.y_axis.set_ticklabels([])
ax.z_axis.set_ticklabels([])
ax.set_xlabel('CLIENT_AGE')
ax.set_ylabel('CLIENT_MMM')
ax.set_title('3D view of K-Means 5 clusters')
ax.dist = 12
plt.show()
```

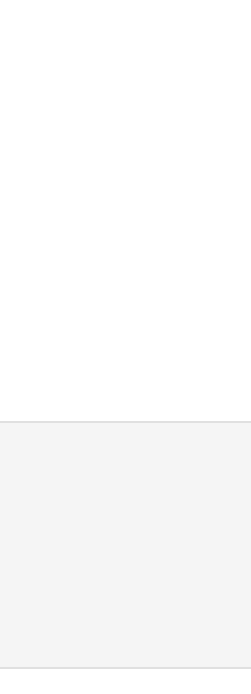


```
trace0 = tracer(KM5_df, 0, 'Cluster 0')
trace1 = tracer(KM5_df, 1, 'Cluster 1')
trace2 = tracer(KM5_df, 2, 'Cluster 2')
trace3 = tracer(KM5_df, 3, 'Cluster 3')
trace4 = tracer(KM5_df, 4, 'Cluster 4')
trace_data5 = [trace0, trace1, trace2, trace3, trace4]
```

```
layout = go.Layout(
    title = 'Clusters wrt Age, Income and Spending Scores and k = 5',
    scene = dict(
        xaxis = dict(title = 'Age'),
        yaxis = dict(title = 'Spending Score'),
        zaxis = dict(title = 'Annual Income')
    )
)
```

```
fig = go.Figure(data=trace_data5, layout=layout)
py.offline.iplot(fig)
```

Clusters wrt Age, Income and Spending Scores and k = 5



### Evaluation metrics

```
db_index = davis_bouldin_score(X, labels6)
print('Daves Bouldin Index: %.2f'%db_index)
so = metrics.silhouette_score(X, labels6)
print('Silhouette Coefficient: %.2f'%so)
ch_score = calinski_harabasz_score(X, labels6)
print('Calinski Harabasz score: %.2f'%ch_score)
```

Daves Bouldin Index: 0.92  
Silhouette Coefficient: 0.42  
Calinski Harabasz score: 4324.80



```
db_index = davis_bouldin_score(X, labels5)
print("Davies Bouldin Index: %0.2f"%db_index)
sc = metrics.silhouette_score(X, labels5)
print("Silhouette Coefficient:%0.2f"%sc)
ch_score = calinski_harabasz_score(X, labels5)
print("Calinski harabasz score:%0.2f"%ch_score)
```

Davies Bouldin Index: 0.84  
Silhouette Coefficient:0.53  
Calinski harabasz score:38863.35

Clearly, the metrics show that using 5 clusters is more optimal

```
# calculate SC for K=2 through K=9
scores = []
for k in range(2,20):
    labels = KMeans(n_clusters=k).fit(X, labels_
    score = metrics.silhouette_score(X, labels_
    scores.append(score)
scores
```

```
[0.4695408092765716,
0.5138604991382025,
0.5122099113572987,
0.5262700225728112,
0.42040109487662564,
0.4359438863341281,
0.4124768885293416,
0.3860510845897715,
0.3848489676446466,
0.36927239836625175,
0.3607123701801566,
0.37087367993795056,
0.349170194419902,
0.37117852087309575,
0.3432624683830541,
0.3542450335211789,
0.356241654849393437,
0.35950497934420367]
```

The highest silhouette score is for k=5 (0.5217648399961166)

## Improve kmeans with PCA

```
# This time the client's profession feature will be added
X_pca = tnp.concatenate([CLIENT_AGE, 'CLIENT_MMM', 'CLIENT_ENCOURS_ENGAGEMENT', 'CLIENT_VRD_MOM', 'TOTAL_PACK', 'CLIENT_PK
data_pca['segment'] = pd.cut(data['CLIENT_AGE'], bins=bins, labels=labels, right=False)
pca = pca_model.fit_transform(X_pca)
pca.shape
```

(50000, 4)

the algorithm automatically selects the best number of principal components that keep 85% of the variance in the original data.

```
# know how many components that the algorithm has selected
pca_model.n_components_
```

4

```
pca_model.explained_variance_
```

```
array([0.24137679, 0.16263663, 0.07410932, 0.05622784])
```

pca explains the variance ratio of 85% with n.component: 4

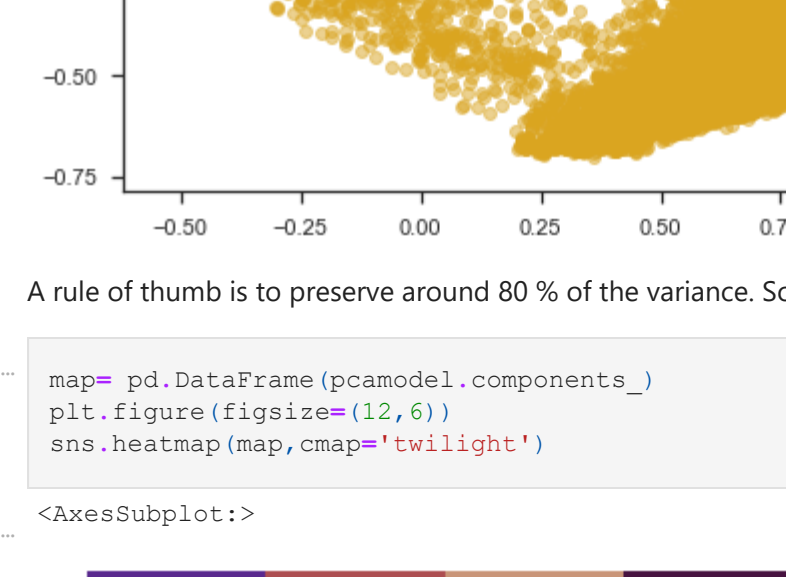
```
pca_model.explained_variance_ratio_
```

```
array([0.40539324, 0.27314883, 0.12446688, 0.09434886])
```



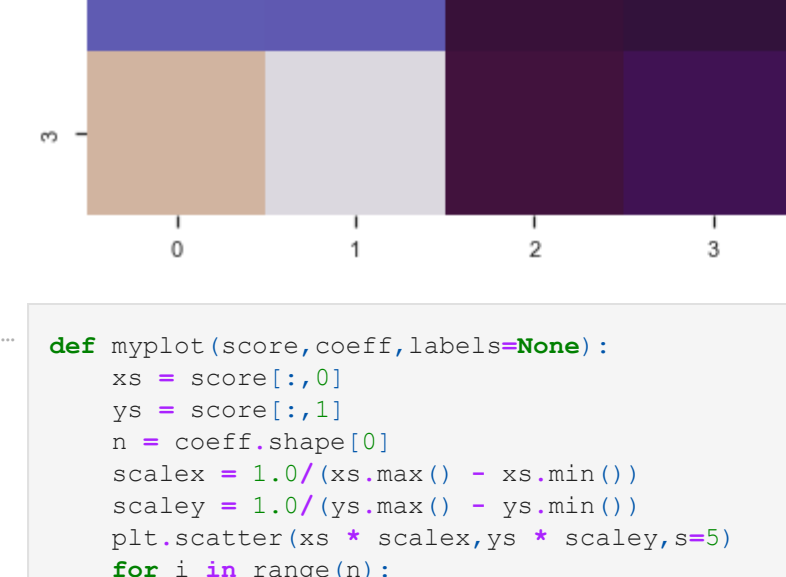
```
plt.bar(range(1, len(pca_model.explained_variance_) + 1), pca_model.explained_variance_)
plt.xlabel('Explained variance')
plt.ylabel('Components')
plt.plot(range(1, len(pca_model.explained_variance_) + 1),
         np.cumsum(pca_model.explained_variance_),
         lw=2,
         label='Cumulative Explained Variance')
plt.legend(loc='upper left')
```

<matplotlib.legend.Legend at 0x211a14af40>



```
plt.plot(pca_model.explained_variance_ratio_)
plt.xlabel('number of components')
plt.ylabel('cumulative explained variance')
plt.show()
```

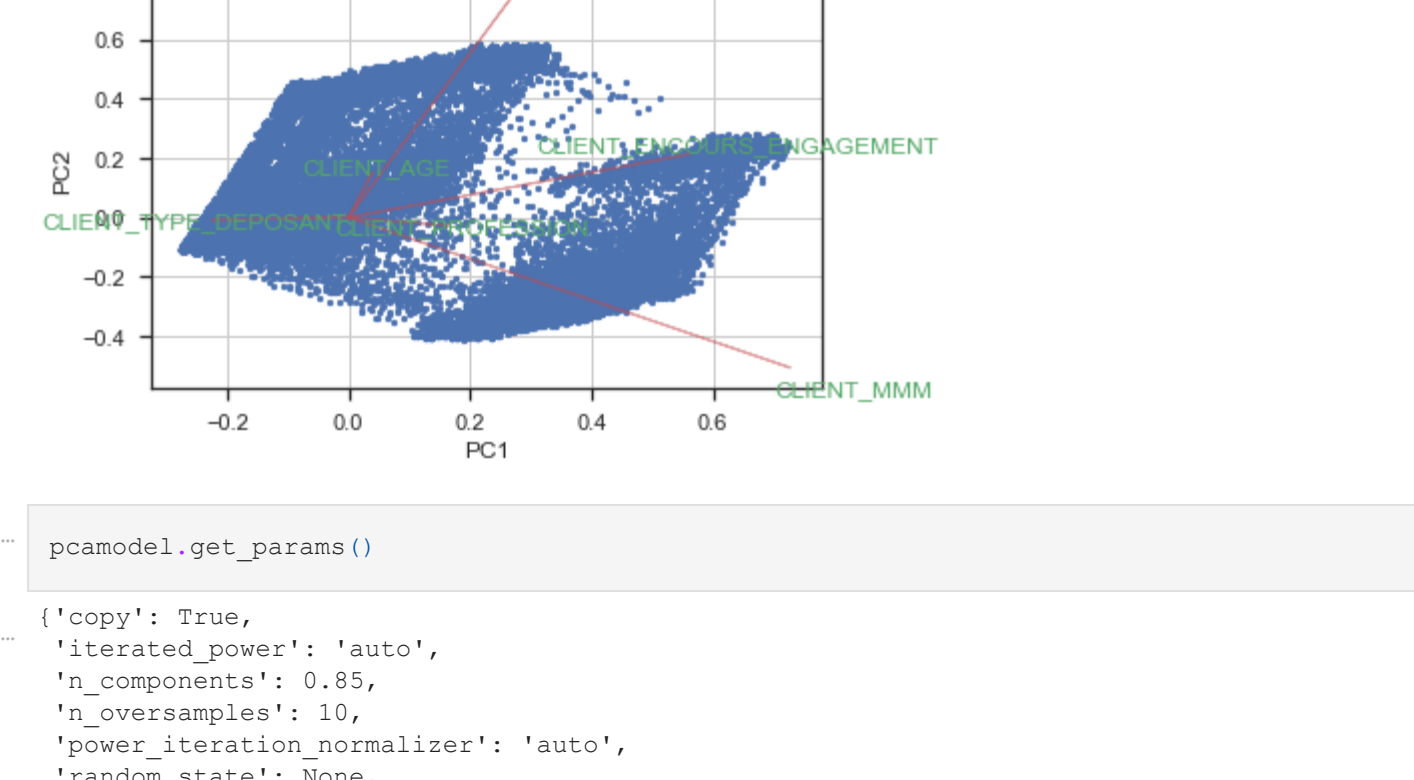
#PCA is at 0 in xscale



A rule of thumb is to preserve around 80 % of the variance. So, in this instance, we decide to keep 3 components.

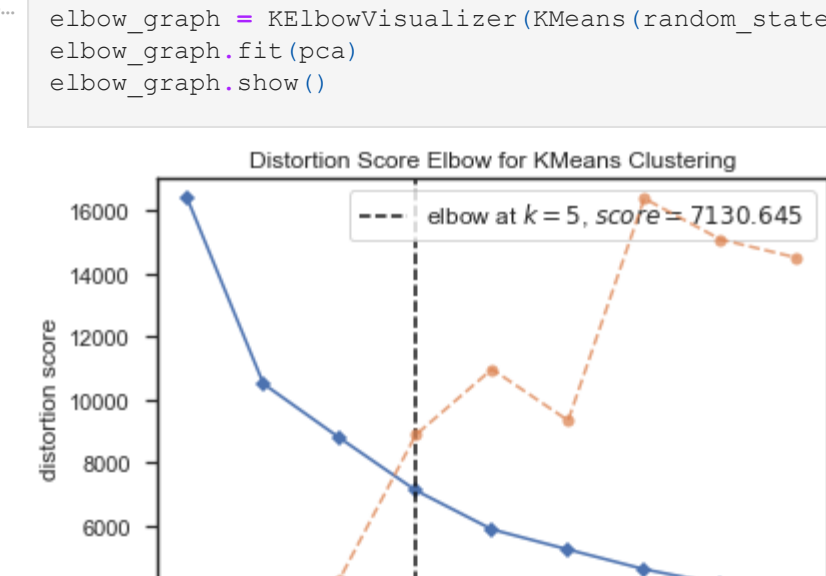
```
map = pd.DataFrame(pca_model.components_)
plt.figure(figsize=(12,6))
sns.heatmap(map, cmap='twilight')
```

<AxesSubplot>



```
def myplot(score,coeff,labels=None):
    xs = score[:,0]
    ys = score[:,1]
    n = coeff.shape[0]
    scalex = 1./((xs.max() - xs.min()))
    scaley = 1./((ys.max() - ys.min()))
    plt.scatter(xs * scalex,ys * scaley,s=5)
    for i in range(n):
        plt.arrow(0, 0, coeff[i,0], coeff[i,1],color = 'r',alpha = 0.5)
        if labels is None:
            plt.text(coeff[i,0]* 1.15, coeff[i,1] * 1.15, "Var"+str(i+1), color = 'green', ha = 'center', va = 'center')
        else:
            plt.text(coeff[i,0]* 1.15, coeff[i,1] * 1.15, labels[i], color = 'g', ha = 'center', va = 'center')
    plt.xlabel("PC1")
    plt.ylabel("PC2")
    plt.grid()

myplot(pca[:,0:2],np.transpose(pca_model.components_[0:2, :]),list(data.columns))
plt.show()
```



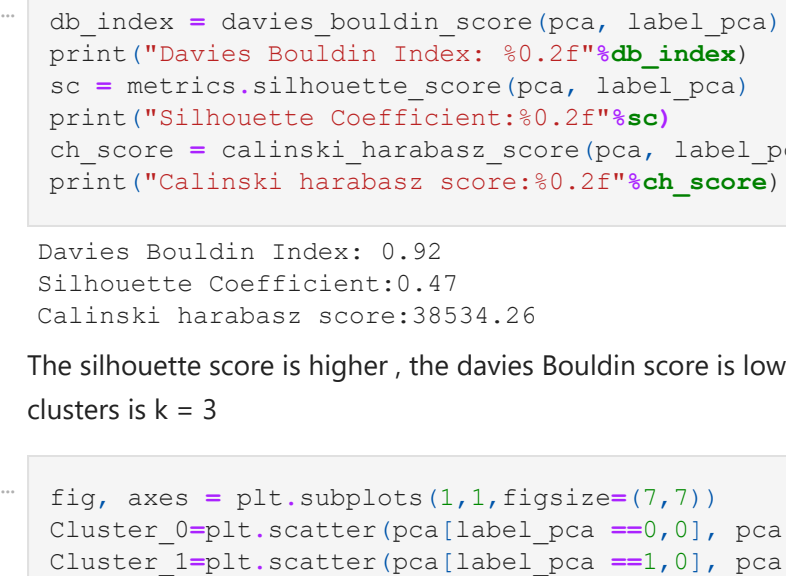
```
CLIENT_NOMBRE_CARTES
CLIENT_ENCOURS_ENGAGEMENT
CLIENT_MMM
CLIENT_VRD_MOM
CLIENT_PROFESSION
```

pca\_model.get\_params()

```
['copy': True,
'iterated_power': 'auto',
'n_components': 10,
'n_oversamples': 10,
'power_iteration_normalizer': 'auto',
'random_state': None,
'svd_solver': 'auto',
'tol': 0.0,
'whiten': False]

wcss=[]
for i in range(1,11):
    kmeans_pca = KMeans(n_clusters= i, init='k-means++', random_state=0)
    kmeans_pca.fit(pca)
    wcss.append(kmeans_pca.inertia_)
```

```
plt.plot(range(1,11), wcss)
plt.title('The Elbow Method')
plt.xlabel('no of clusters')
plt.ylabel('wcss')
plt.show()
```



elbow\_graph = KElbowVisualizer(KMeans(random\_state=43), k=10)

elbow\_graph.fit(pca)

elbow\_graph.show()



<AxesSubplot:title='center':Distortion Score Elbow for KMeans Clustering', xlabel='k', ylabel='distortion score'>

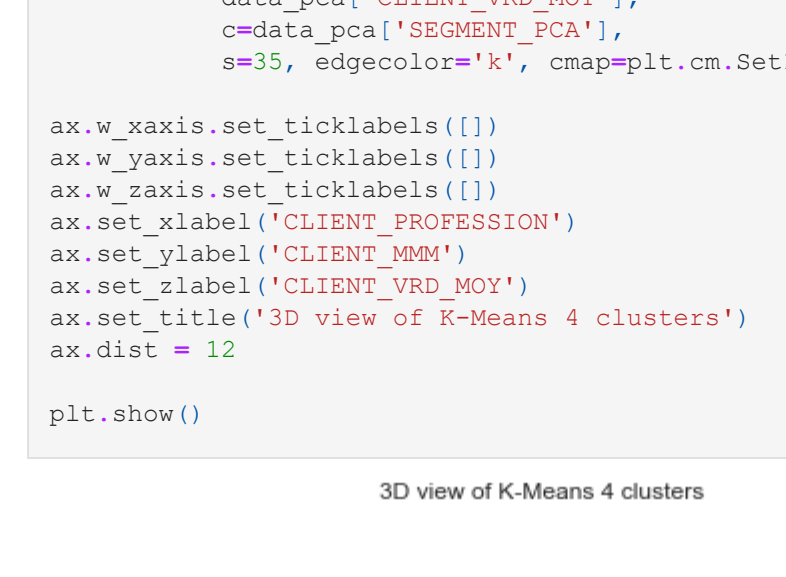
```
kmeans_pca = KMeans(n_clusters=4, init='k-means++', random_state=42)
```

label\_pca = kmeans\_pca.fit\_predict(pca)

```
db_index = davis_bouldin_score(pca, label_pca)
print("Davies Bouldin Index: %0.2f"%db_index)
sc = metrics.silhouette_score(pca, label_pca)
print("Silhouette Coefficient:%0.2f"%sc)
ch_score = calinski_harabasz_score(pca, label_pca)
print("Calinski harabasz score:%0.2f"%ch_score)
```

Davies Bouldin Index: 0.96  
Silhouette Coefficient:0.38  
Calinski harabasz score:33961.70

```
plt.scatter(pca[label_pca == 0, 0], pca[label_pca == 0, 1], s=100, c='red', label = 'Cluster 1')
plt.scatter(pca[label_pca == 1, 0], pca[label_pca == 1, 1], s=100, c='blue', label = 'Cluster 2')
plt.scatter(pca[label_pca == 2, 0], pca[label_pca == 2, 1], s=100, c='green', label = 'Cluster 3')
plt.scatter(pca[label_pca == 3, 0], pca[label_pca == 3, 1], s=100, c='cyan', label = 'Cluster 4')
plt.legend()
plt.title('Clusters of customers')
plt.show()
```



kmeans\_pca = KMeans(n\_clusters=3, init='k-means++', random\_state=42)

label\_pca = kmeans\_pca.fit\_predict(pca)

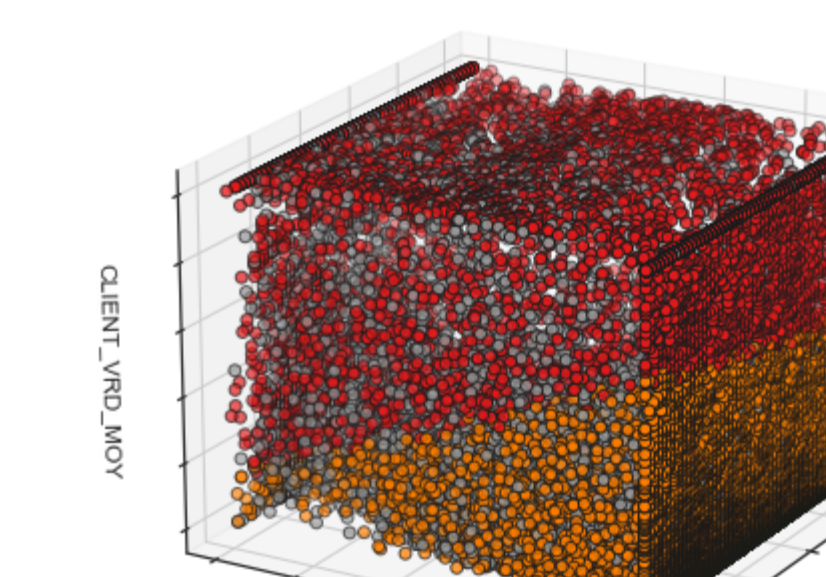
```
db_index = davis_bouldin_score(pca, label_pca)
print("Davies Bouldin Index: %0.2f"%db_index)
sc = metrics.silhouette_score(pca, label_pca)
print("Silhouette Coefficient:%0.2f"%sc)
ch_score = calinski_harabasz_score(pca, label_pca)
print("Calinski harabasz score:%0.2f"%ch_score)
```

Davies Bouldin Index: 0.92  
Silhouette Coefficient:0.47  
Calinski harabasz score:38334.26

The silhouette score is higher, the Davies Bouldin score is lower, and the calinski harabasz is higher so the chosen value for the number of clusters is k = 3

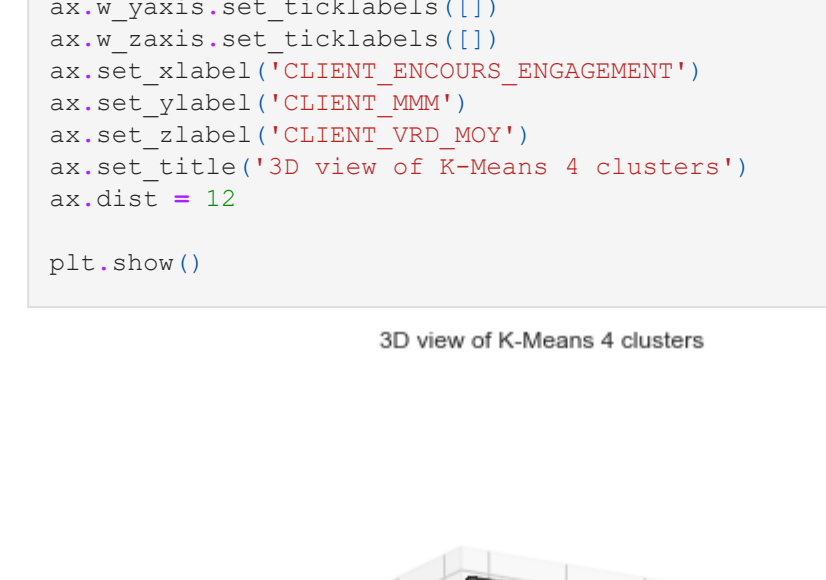
```
fig, axes = plt.subplots(1,1,figsize=(7,7))
ax = Axes3D(fig, rect=[0, 0, .99, 1], elev=20, azim=210)
ax.scatter(data_pca['CLIENT_ENCOURS_ENGAGEMENT'],
           data_pca['CLIENT_MMM'],
           data_pca['CLIENT_VRD_MOM'],
           data_pca['SEGMENT_PCA'],
           cdata_pca['SEGMENT_PCA'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)

ax.w_xaxis.set_ticklabels([])
ax.w_yaxis.set_ticklabels([])
ax.w_zaxis.set_ticklabels([])
ax.set_xlabel('CLIENT_ENCOURS_ENGAGEMENT')
ax.set_ylabel('CLIENT_MMM')
ax.set_zlabel('CLIENT_VRD_MOM')
ax.set_title('3D view of K-Means 4 clusters')
ax.dist = 12
plt.show()
```



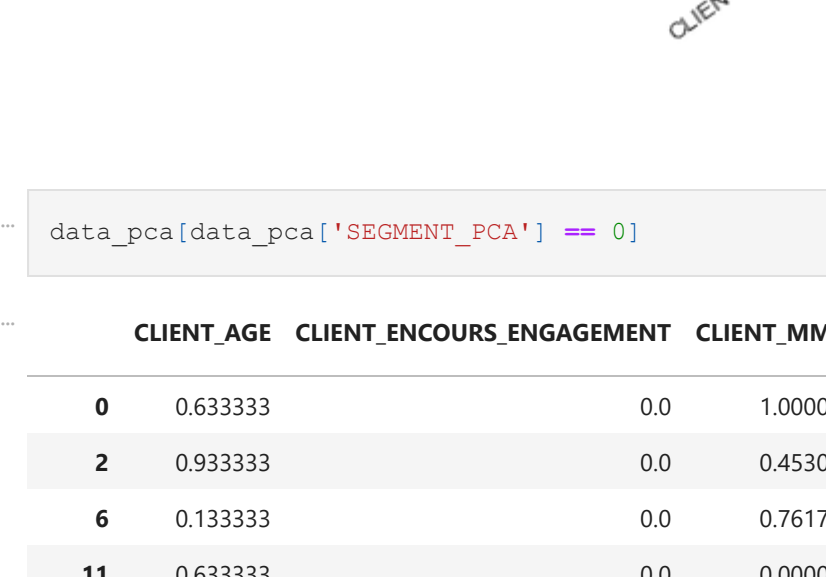
```
fig = plt.figure(figsize=(7, 7))
ax = Axes3D(fig, rect=[0, 0, .99, 1], elev=20, azim=210)
ax.scatter(data_pca['CLIENT_ENCOURS_ENGAGEMENT'],
           data_pca['CLIENT_MMM'],
           data_pca['CLIENT_VRD_MOM'],
           data_pca['SEGMENT_PCA'],
           cdata_pca['SEGMENT_PCA'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)

ax.w_xaxis.set_ticklabels([])
ax.w_yaxis.set_ticklabels([])
ax.w_zaxis.set_ticklabels([])
ax.set_xlabel('CLIENT_ENCOURS_ENGAGEMENT')
ax.set_ylabel('CLIENT_MMM')
ax.set_zlabel('CLIENT_VRD_MOM')
ax.set_title('3D view of K-Means 3 clusters')
ax.dist = 12
plt.show()
```



```
fig = plt.figure(figsize=(7, 7))
ax = Axes3D(fig, rect=[0, 0, .99, 1], elev=20, azim=210)
ax.scatter(data_pca['CLIENT_ENCOURS_ENGAGEMENT'],
           data_pca['CLIENT_MMM'],
           data_pca['CLIENT_VRD_MOM'],
           data_pca['SEGMENT_PCA'],
           cdata_pca['SEGMENT_PCA'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)

ax.w_xaxis.set_ticklabels([])
ax.w_yaxis.set_ticklabels([])
ax.w_zaxis.set_ticklabels([])
ax.set_xlabel('CLIENT_ENCOURS_ENGAGEMENT')
ax.set_ylabel('CLIENT_MMM')
ax.set_zlabel('CLIENT_VRD_MOM')
ax.set_title('3D view of K-Means 4 clusters')
ax.dist = 12
plt.show()
```



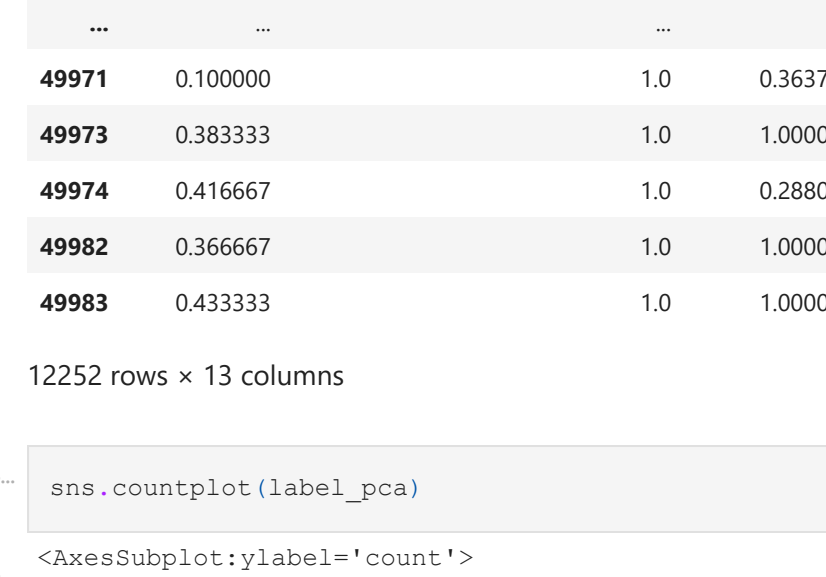
```
data_pca[data_pca['SEGMENT_PCA'] == 0]
```

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	C
0	0.633333	0.000000	1.000000	0.6	0.258065	0.00	
1	0.433333	0.145041	0.156169	0.2	0.645161	0.00	
2	0.933333	0.000000	0.430350	0.2	0.000000	0.00	
3	0.500000	0.000000	0.000000	0.0	0.580645	0.25	
4	0.500000	1.000000	1.000000	0.4	0.258065	0.00	
...	...	...	...	...	...	...	...
49995	0.466667	0.000000	0.000000	0.0	0.064516	0.25	
49996	0.200000	0.000000	0.000000	0.2	0.645161	0.25	
49997	0.366667	0.000000	0.000000	0.2	0.645161	0.25	
49998	0.333333	0.000000	0.000000	0.0	0.645161	0.25	
49999	0.233333	0.000000	0.000000	0.0	0.967742	0.25	

50000 rows × 13 columns

```
fig = plt.figure(figsize=(7, 7))
ax = Axes3D(fig, rect=[0, 0, .99, 1], elev=20, azim=210)
ax.scatter(data_pca['CLIENT_ENCOURS_ENGAGEMENT'],
           data_pca['CLIENT_MMM'],
           data_pca['CLIENT_VRD_MOM'],
           data_pca['SEGMENT_PCA'],
           cdata_pca['SEGMENT_PCA'],
           s=35, edgecolor='k', cmap=plt.cm.Set1)

ax.w_xaxis.set_ticklabels([])
ax.w_yaxis.set_ticklabels([])
ax.w_zaxis.set_ticklabels([])
ax.set_xlabel('CLIENT_ENCOURS_ENGAGEMENT')
ax.set_ylabel('CLIENT_MMM')
ax.set_zlabel('CLIENT_VRD_MOM')
ax.set_title('3D view of K-Means 4 clusters')
ax.dist = 12
plt.show()
```



```
data_pca[data_pca['SEGMENT_PCA'] == 1]
```

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	C
1	0.433333	0.145041	0.156169	0.2	0.645161	0.00	
2	0.500000	0.000000	0.000000	0.4	0.258065	0.00	
3	0.933333	0.000000	0.430350	0.0	0.580645	0.25	
5	0.333333	0.000000	0.000000	0.0	0.387097	0.25	
6	0.633333	0.000000	0.000000	0.0	0.193548	0.25	
11	0.733333	0.000000	0.000000	0.0	0.838710	0.25	
17	0.266667	0.000000	0.000000	0.0	0.838710	0.25	
...	...	...	...	...	...	...	...
49994	0.383333	0.000000	0.998480	0.6	0.258065	0.00	
49990	0.566667	0.000000	0.312159	0.2	0.000000	0.00	
49992	0.366667	0.000000	0.000000	0.0	0.838710	0.25	
49993	0.366667	0.000000	0.000000	0.2	0.645161	0.25	
49996	0.200000	0.000000	0.000000	0.2	0.645161	0.25	

9572 rows × 13 columns

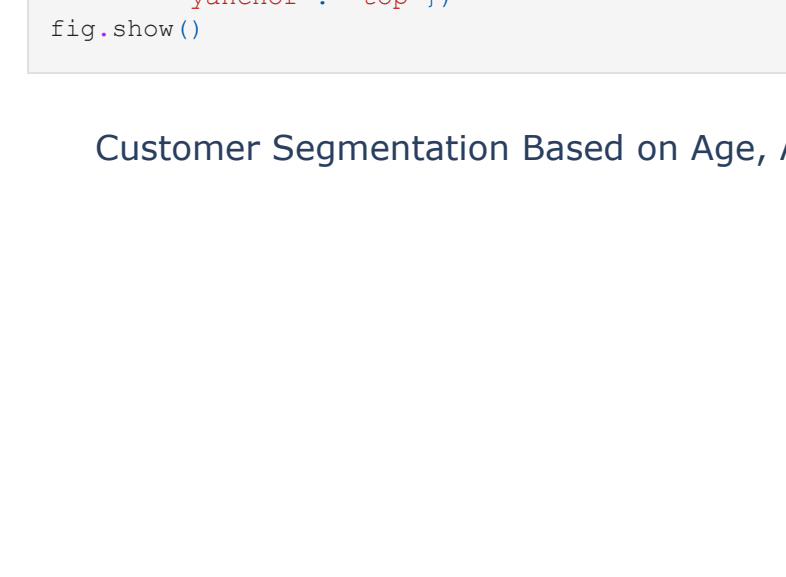
```
data_pca[data_pca['SEGMENT_PCA'] == 2]
```

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPOSANT	C
1	0.433333	0.145041	0.156169	0.2	0.645161	0.00	
7	0.400000	1.0	1.081134	0.4	0.258065	0.00	
9	0.266667	1.0	0.262139	0.0	0.000000	0.00	
10	0.566667	1.0	0.190365	0.6	0.645161	0.25	
13	0.366667	1.0	0.000000	0.6	0.161290	0.00	
...	...	...	...	...	...	...	...
49994	0.100000	1.0	0.363748	0.0	0.645161	0.00	
49973	0.383333	1.0	1.000000	0.8	0.258065	0.00	
49974	0.416667	1.0	0.288073	0.0	0.580645	0.25	
49982	0.366667	1.0	1.000000	0.6	0.290323	0.00	
49983	0.433333	1.0	1.000000	0.6	0.451613	0.00	

12252 rows × 13 columns

sns.countplot(label\_pca)

<AxesSubplot:ylabel='count'>



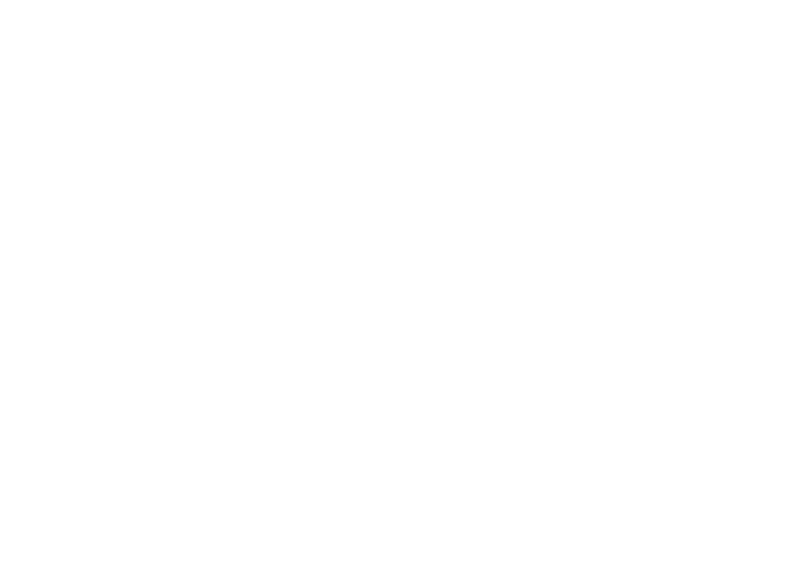
```
data_pca.groupby('SEGMENT_PCA').mean().round(3)
```

	CLIENT_AGE	CLIENT_ENCOURS_ENGAGEMENT	CLIENT_MMM	CLIENT_NOMBRE_CARTES	CLIENT_PROFESSION	CLIENT_TYPE_DEPO
SEGMENT_PCA						
0	0.491	0.006	0.405	0.183	0.524	
1	0.371	0.012	0.053	0.097	0.596	
2	0.986	0.989	0.492	0.349	0.445	

data\_pca['SEGMENT\_PCA'].value\_counts()

```
1    28176
2    12252
0     9572
Name: SEGMENT_PCA, dtype: int64
```

```
labels = ['0', '1', '2']
size = [28052, 9614, 9522]
colors = ['lightgreen', 'orange', 'cyan', 'pink']
explode = None
plt.rcParams['figure.figsize'] = (7, 7)
plt.pie(size, colors=colors, explode=explode, labels=labels, shadow=True, autopct='%2.2f%%')
plt.title('A pie chart Representing the count of each segment')
plt.axis('off')
plt.legend()
plt.show()
```



```
fig=plt.scatter_3d(data_frame=X_pca,x="CLIENT_AGE",y="CLIENT_MMM",z="CLIENT_VRD_MOM",color=label_pca, color_cont
data_pca['segment'] = pd.cut(data['CLIENT_AGE'], bins=bins, labels=labels, right=False)
print (data_pca)
```

Customer Segmentation Based on Age, Average monthly flow and Actual values of deposits



```
map = pd.DataFrame(pca_model.components_)
plt.figure(figsize=(12,6))
sns.heatmap(map, cmap='twilight')
```

<AxesSubplot>



```
bins = [30,40,50,70,80]
labels = ['Young Adults', 'Adults', 'Middle Aged', 'Senior Citizens']
data_pca['segment'] = pd.cut(data['CLIENT_AGE'], bins=bins, labels=labels, right=False)
print (data_pca)
```



