

K-Medoids clustering on socio-economic dataset

```
In [6]: import pandas as pd
from sklearn_extra.cluster import KMedoids
from sklearn.metrics import silhouette_score

df = pd.read_excel('Preprocessed_data_standardscaler.xlsx')

# Dropping the ID Column, as it is not required in clustering.

new_df = df.drop(columns=['ID'])
new_df
```

	Normalised Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)	Experience (>30yr)	Helmet (Full-mask)
0	-0.927352	1	1	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0
1	0.387063	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
2	-0.739578	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
3	0.762611	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
4	-1.678446	1	0	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
476	-0.739578	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
477	-1.302899	1	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0
478	2.077026	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	1	0
479	-0.551804	1	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0
480	-0.551804	1	0	0	0	1	0	0	1	0	0	0	0	0	0	0	1	0	0

481 rows x 19 columns

```
In [7]: # Now the Silhouette Score method is used to determine the number of clusters to achieve the best possible clustering
# Here the scores are calculated for k = {2,3,4,5,6}

scores = {}

for i in range(2,7):
    kmedoids = KMedoids(n_clusters=i,method='pam',init='build',random_state=42)
    kmedoids.fit(new_df)
    score = silhouette_score(new_df,kmedoids.labels_)
    scores[i]=score

df_scores = pd.DataFrame(list(scores.items()), columns=['no. of clusters', 'silhouette score'])
df_scores

# since setting k = 5 yields the highest score, the number of clusters are chosen to be 5.
# As can be seen in the table below:
```

	no. of clusters	silhouette score
0	2	0.137599
1	3	0.152380
2	4	0.170523
3	5	0.172540
4	6	0.162440

```
In [8]: # Now applying the K-Medoids algorithm.

optimal_clusters = 5

kmedoids_1 = KMedoids(n_clusters=optimal_clusters,method='pam',init='build',random_state=42)
kmedoids_1.fit(new_df)

list1 = list(kmedoids_1.labels_)

df_clustered = new_df.copy()

df_clustered['Cluster'] = list1
df_clustered

# As can be seen in the below dataset, a new 'Cluster' column has been added at the very end:
```

	Normalised Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)	Experience (>30yr)	Helmet (Full-mask)
0	-0.927352	1	1	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0
1	0.387063	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
2	-0.739578	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
3	0.762611	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
4	-1.678446	1	0	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
476	-0.739578	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
477	-1.302899	1	0	0	1	0	0	0	0	1	0	0	0	0	1	0	0	0	0
478	2.077026	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	0	1	0
479	-0.551804	1	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0
480	-0.551804	1	0	0	0	1	0	0	1	0	0	0	0	0	0	0	1	0	0

481 rows x 20 columns

```
In [27]: final_df = df_clustered.copy()

old_df = pd.read_excel('Asansol socio-economic data 1.xlsx')
age = list(old_df['Age (Yrs)'])

final_df['Age (Yrs)'] = age

final_df = final_df.drop(columns = ['Normalised Age (Yrs)'])

list2 = list(final_df.columns)
new_order = [list2[-1]] + list2[0:-1]

final_df = final_df[new_order]
final_df
```

	Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)	Experience (>30yr)	Helmet (Full-mask)
0	26	1	1	0	0	0	1	0	0	0	0	0	0	0	0	1	0	0	1
1	40	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0
2	28	1	0	0	1	0	0	0	1	0	0	0	0	0	1	0	0	0	1
3	44	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	0	0	0
4	18	1	0	1	0	0	0	1	0	0	0	1	0	0	0	0	0	0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
476	28	1	0	0	1	0	0	0	1	0	0	0	0	0	1	0	0	0	1
477	22	1	0	0	1	0	0	0	0	1	0	0	0	1	0	0	0	0	1
478	58	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
479	30	1	0	0	1	0	0	0	1	0	0	0	0	0	0	1	0	0	1
480	30	1	0	0	0	1	0	0	1	0	0	0	0	0	0	1	0	0	1

481 rows x 20 columns

```
In [32]: cluster_1 = final_df.loc[final_df['Cluster']==0]
cluster_1
cluster_1.describe()
```

	Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)	Experience (>30yr)
count	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000	71.000000
mean	51.873239	0.943662	0.140845	0.183099	0.577465	0.098592	0.070423	0.239437	0.267606	0.380282	0.042254	0.0	0.014085	0.028169	0.056338	0.169014	0.591549	0.0
std	8.349388	0.232214	0.350338	0.389500	0.497479	0.300235	0.257679	0.429777	0.445862	0.488911	0.202599	0.0	0.118678	0.166633	0.232214	0.377432	0.495046	0.0
min	35.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	46.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	51.000000	1.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000
75%	58.500000	1.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	1.000000	1.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000
max	73.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

```
In [30]: cluster_2 = final_df.loc[final_df['Cluster']==1]
cluster_2.describe()
```

	Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)	Experience (>30yr)
count	82.000000	82.000000	82.0	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000	82.0	82.000000	82.000000	82.000000	82.000000	82.000000	82.000000
mean	38.329268	0.951220	0.0	0.085366	0.878049	0.036585	0.073171	0.073171	0.658537	0.170732	0.024390	0.0	0.048780	0.048780	0.024390	0.853659	0.024390	0.0
std	6.299235	0.216734	0.0	0.281145	0.329243	0.188897	0.262019	0.262019	0.477119	0.378590	0.155207	0.0	0.216734	0.216734	0.155207	0.355623	0.155207	0.0
min	28.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	35.000000	1.000000	0.0	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000
50%	37.000000	1.000000	0.0	0.000000	1.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000
75%	40.750000	1.000000	0.0	0.000000	1.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000
max	58.000000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	0.0	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000

```
In [31]: cluster_3 = final_df.loc[final_df['Cluster']==2]
cluster_3.describe()
```

Out [31]:

	Age (Yrs)	Gender (Male)	Education (<Class 10)	Education (Class 10-12)	Graduate	Higher Edu	Income <20k	Income 20-40k	Income 40-60k	Income 60-80k	Income >80k	Experience (<1yr)	Experience (1-2yr)	Experience (2-5yr)	Experience (5-10yr)	Experience (10-20yr)	Experience (20-30yr)
count	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000	97.0	97.000000	97.000000	97.000000	97.000000	97.000000	97.000000
mean	32.298969	0.958763	0.742268	0.195876	0.030928	0.030928	0.865979	0.092784	0.020619	0.020619	0.0	0.051546	0.103093	0.14433	0.494845	0.175258	0.030928
std	8.214780	0.199871	0.439658	0.398935	0.174022	0.174022	0.342444	0.291636	0.142842	0.142842	0.0	0.222258	0.305660	0.35325	0.502571	0.382162	0.174022
min	18.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.00000	0.000000	0.000000	0.000000
25%	26.000000	1.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.00000	0.000000	0.000000	0.000000
50%	32.000000	1.000000	1.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.00000	0.000000	0.000000	0.000000
75%	38.000000	1.000000	1.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000	0.0	0.000000	0.000000	0.00000	1.000000	0.000000	0.000000
max	50.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	1.000000	0.0	1.000000	1.000000	1.00000	1.000000	1.000000	1.000000