

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn import preprocessing
import seaborn as sns
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler

data = pd.read_csv("mall_customers.csv", sep=";")
data
```

Out[2]:

CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15
1	2	Male	21	15
2	3	Female	20	16
3	4	Female	23	16
4	5	Female	31	17
...	...	...	...	...
195	196	Female	35	120
196	197	Female	45	126
197	198	Male	32	126
198	199	Male	32	127
199	200	Male	30	137
...	...	...	...	...
200	rows = 5	columns = 5		

ANALIZA

Stworzyć histogramy dla wieku, rocznego dochodu i oceny wydawania.

```
In [13]: plt.hist(data["Age"], facecolor='blue', alpha=0.5)
plt.title("Spending Score Histogram")
plt.xlabel("Age in Years")
plt.ylabel("Amount of People")
plt.show()

In [17]: plt.hist(data["Annual Income (k$)"], facecolor='blue', alpha=0.5)
plt.title("Annual Income Histogram")
plt.xlabel("Income in k $")
plt.ylabel("Amount of People")
plt.show()

In [18]: plt.hist(data["Spending Score (1-100)"], facecolor='blue', alpha=0.5)
plt.title("Spending Score Histogram")
plt.xlabel("Score from 0 to 100")
plt.ylabel("Amount of People")
plt.show()
```

Porównać na wykresie liczbę kobiet i mężczyzn.

```
In [14]: amount_of_women = len(list(filter(lambda x: x == "Female", data["gender"])))
amount_of_men = len(list(filter(lambda x: x == "Male", data["gender"])))
amount = len(data["gender"])
# Now every iteration we will have to more by one because of men and women
objects = ("Men " + str(amount_of_men), "Women " + str(amount_of_women))
x_pos = np.arange(len(objects))
performance = [amount_of_men, amount_of_women, amount]

plt.bar(x_pos, performance, align='center', alpha=0.5)
plt.xticks(x_pos, objects)
plt.xlabel('Amount')
plt.title('Amount of Women and Men')
plt.show()
```

Przedstawić na wykresie zależności pomiędzy wiekiem, rocznym dochodem i oceną wydawania (6 wykresów np. wiek i dochody, wiek i ocena wydawania itd.).

```
In [19]: titleDict = {'fontsize': 28, 'weight': 'bold'}
sns.set_context('paper')
sns.relplot(data=data,
            x="Age",
            y="Spending Score (1-100)",
            aspect=2.5,
            kind='line')
plt.title("Age to Annual Income", fontdict=titleDict)
plt.show()
```

Age to Annual Income

```
In [12]: sns.set_context('paper')
sns.relplot(data=data,
            x="Spending Score (1-100)",
            y="Age",
            aspect=2.5,
            kind='line')
plt.title("Age to Score", fontdict=titleDict)
plt.show()
```

Age to Score

```
In [13]: sns.set_context('paper')
sns.relplot(data=data,
            x="Annual Income (k$)",
            y="Spending Score (1-100)",
            aspect=2.5,
            kind='line')
plt.title("Annual Income to Score", fontdict=titleDict)
plt.show()
```

Annual Income to Score

```
In [14]: sns.set_context('paper')
sns.relplot(data=data,
            x="Annual Income (k$)",
            y="Age",
            aspect=2.5,
            kind='line')
plt.title("Annual Income to Age", fontdict=titleDict)
plt.show()
```

Annual Income to Age

```
In [15]: sns.set_context('paper')
sns.relplot(data=data,
            x="Spending Score (1-100)",
            y="Age",
            aspect=2.5,
            kind='line')
plt.title("Score to Age", fontdict=titleDict)
plt.show()
```

Score to Age

```
In [16]: sns.set_context('paper')
sns.relplot(data=data,
            x="Spending Score (1-100)",
            y="Annual Income (k$)",
            aspect=2.5,
            kind='line')
plt.title("Score to Annual Income", fontdict=titleDict)
plt.show()
```

Score to Annual Income

Dokonać niezbędnych zmian w bazie danych potrzebnych do rozpoczęcia procesu klasteryzacji (kolumna płeć i ID)

```
In [6]: ie = preprocessing.LabelEncoder()
new_data = ie.fit_transform(data["gender"])
data["gender"] = new_data
data.drop(columns=["CustomerID"], inplace=True)
```

Wygenerować macierz kowariancji i korelacji.

```
In [7]: correlation = data.corr()
correlation.round(3).style.background_gradient()
```

Out[7]:

	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
Gender	1.000000	0.291000	0.056000	-0.059000
Age	0.291000	1.000000	-0.022000	-0.327000
Annual Income (k\$)	0.056000	-0.022000	1.000000	0.020000
Spending Score (1-100)	-0.059000	-0.327000	0.020000	1.000000

```
In [8]: covariance = data.cov()
covariance.round(3).style.background_gradient()
```

Out[8]:

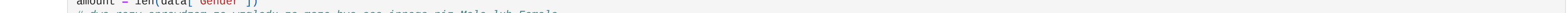
	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
Gender	0.148000	0.423000	0.072000	-0.747000
Age	0.423000	0.282000	-0.009000	-1.161000
Annual Income (k\$)	0.072000	-0.009000	1.601000	0.713000
Spending Score (1-100)	-0.747000	-1.161000	0.713000	69.616000

KLASTERYZACJA PODSTAWOWA

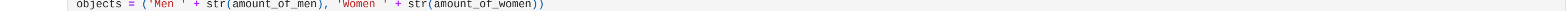
Dokonać klasteryzacji z użyciem algorytmu k-średnich dla k=5 z uwzględnieniem tylko dwóch kolumn (wybrać je na podstawie wyników macierzy kowariancji i korelacji).

```
In [9]: # KMeans clustering
# Age and Spending Score are the features we will use to cluster the data
cluster_without_PCA_data = data.iloc[:, 1:4].drop(columns=["Annual Income (k$)"], values)
kmeans = KMeans(n_clusters=5, random_state=0)
kmeans.fit_transform(cluster_without_PCA_data)
kmeans.fit_transform(cluster_without_PCA_data)
```

```
In [10]: plt.scatter(cluster_without_PCA_data[:, 0], cluster_without_PCA_data[:, 1], c=kmeans.labels_, cmap='jet')
plt.xlabel("Age")
plt.ylabel("Spending Score (1-100)")
plt.show()
```



```
In [11]: plt.scatter(kmeans.cluster_centers_[0], kmeans.cluster_centers_[1], s=200, c='yellow', label = 'Centroids')
plt.title('Centers of clusters')
plt.xlabel("Age")
plt.ylabel("Spending Score (1-100)")
plt.show()
```



KLASTERYZACJA Z PCA

Która kolumna/kolumny przechowują najwięcej informacji ?

Jak można zauważyć w data.cov() wyżej największą informację przechowuje Annual Income (k\$) oraz Spending Score (1-100)

Czy możemy zredukować do dwóch wymiarów przy pozostaniu na poziomie powyżej 80% wariancji?

```
In [20]: features = ["Age", "Annual Income (k$)", "Spending Score (1-100)"]
# Separating out the features
x = data.iloc[:, features]
# Standardizing the features
x = StandardScaler().fit_transform(x)
pca = PCA(n_components=2)
principalComponents = pca.fit_transform(x)
principalDF = pd.DataFrame(data = principalComponents)
principalDF.columns = ["principal component 1", "principal component 2"]
pca.explained_variance_ratio_.sum()
```

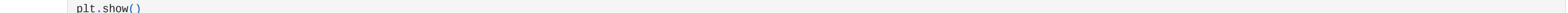
Out[20]: 0.77545466976718

Nie można ponieważ nie wychodzi poza 80 procent

```
In [13]: pca = PCA(n_components=3)
principalComponents = pca.fit_transform(x)
principalDF = pd.DataFrame(data = principalComponents)
print(pca.explained_variance_ratio_.sum())
```

Zawiera powyżej 80 procent więc działamy na tych danych

```
In [23]: kmeans = KMeans(n_clusters=5, random_state=0)
kmeans.fit_transform(x)
fig = plt.figure(figsize=(8, 8))
ax = fig.add_subplot(1, 1, projection='3d')
ax.scatter(principalDF[0], principalDF[1], principalDF[2], c=kmeans.labels_, cmap='jet', s=10)
ax.set_xlabel('Principal Component 1')
ax.set_ylabel('Principal Component 2')
ax.set_zlabel('Principal Component 3')
ax.grid()
```

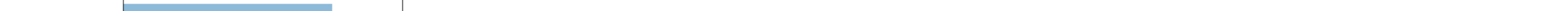


ANALIZA

Porównać wyniki klasteryzacji z użyciem PCA i bez

```
In [47]: plt.scatter(cluster_without_PCA_data[:, 0], cluster_without_PCA_data[:, 1], c=kmeans.labels_, cmap='jet')
plt.title("Cluster without PCA")
plt.xlabel("Age")
plt.ylabel("Spending Score (1-100)")
plt.show()
```

```
In [72]: plt.scatter(principalDF[0], principalDF[1], c=kmeans.labels_, cmap='jet', s=10)
plt.title("Cluster with PCA")
plt.xlabel("Principal Component 1")
plt.ylabel("Principal Component 2")
plt.show()
```



Klasteryzacja z PCA jest dużo wydajniejsza i dokładna od klasteryzacji bez PCA

Opisać krótko każdą grupę nadając jej odpowiednią nazwę własną (np. Klienci docelowi). Ocenić potencjał grup pod względem przyszłych wyników sprzedaży. Do kogo należy kierować najwięcej informacji marketingowych?

```
In [106]: final_data = pd.concat([data, pd.DataFrame(data = kmeans.predictions_)], axis=1)
grouped_final_data = final_data.groupby(["Group_Predicted"])
print(grouped_final_data.mean().round())
print(grouped_final_data.get_group(0).mean().round())
print(grouped_final_data.get_group(1).mean().round())
print(grouped_final_data.get_group(2).mean().round())
print(grouped_final_data.get_group(3).mean().round())
print(grouped_final_data.get_group(4).mean().round())
```

Gender

	Gender	Age	Annual Income (k\$)	Spending Score (1-100)	Group_Predicted
0	0.0	46.0	88.0	19.0	0
1	0.0	46.0	88.0	19.0	0
2	0.0	46.0	88.0	19.0	0
3	0.0	46.0	88.0	19.0	0
4	0.0	46.0	88.0	19.0	0
5	0.0	46.0	88.0	19.0	0
6	0.0	46.0	88.0	19.0	0
7	0.0	46.0	88.0	19.0	0
8	0.0	46.0	88.0	19.0	0
9	0.0	46.0	88.0	19.0	0
10	0.0	46.0	88.0	19.0	0
11	0.0	46.0	88.0	19.0	0
12	0.0	46.0	88.0	19.0	0
13	0.0	46.0	88.0	19.0	0
14	0.0	46.0	88.0	19.0	0
15	0.0	46.0	88.0	19.0	0
16	0.0	46.0	88.0	19.0	0
17	0.0	46.0	88.0	19.0	0
18	0.0	46.0	88.0	19.0	0
19	0.0	46.0	88.0	19.0	0
20	0.0	46.0	88.0	19.0	0
21	0.0	46.0	88.0	19.0	0
22	0.0	46.0	88.0	19.0	0
23	0.0	46.0	88.0	19.0	0
24	0.0	46.0	88.0	19.0	0
25	0.0	46.0	88.0	19.0	0
26	0.0	46.0	88.0	19.0	0
27	0.0	46.0	88.0	19.0	0
28	0.0	46.0	88.0	19.0	0
29	0.0	46.0	88.0	19.0	0
30	0.0	46.0	88.0	19.0	0
31	0.0	46.0	88.0	19.0	0
32	0.0	46.0	88.0	19.0	0
33	0.0	46.0	88.0	19.0	0
34	0.0	46.0	88.0	19.0	0
35	0.0	46.0	88.0	19.0	0
36	0.0	46.0	88.0	19.0	0
37	0.0	46.0	88.0	19.0	0
38	0.0	46.0	88.0	19.0	0
39	0.0	46.0	88.0	19.0	0
40	0.0	46.0	88.0	19.0	0
41	0.0	46.0	88.0	19.0	0
42	0.0	46.0	88.0	19.0	0
43	0.0	46.0	88.0	19.0	0
44	0.0	46.0	88.0	19.0	0
45	0.0	46.0	88.0	19.0	0
46	0.0	46.0	88.0	19.0	0
47	0.0	46.0	88.0	19.0	0
48	0.0	46.0	88.0	19.0	0
49	0.0	46.0	88.0	19.0	0
50	0.0	46.0	88.0	19.0	0
51	0.0	46.0	88.0	19.0	0
52	0.0	46.0	88.0	19.0	0
53	0.0	46.0	88.0	19.0	0
54	0.0	46.0	88.0	19.0	0
55	0.0	46.0	88.0	19.0	0
56	0.0	46.0	88.0	19.0	0
57	0.0	46.0	88.0	19.0	0
58	0.0	46.0	88.0	19.0	0
59	0.0	46.0	88.0	19.0	0
60	0.0	46.0	88.0	19.0	0
61	0.0	46.0	88.0	19.0	0
62	0.0	46.0	88.0	19.0	0
63	0.0	46.0	88.0	19.0	0
64	0.0	46.0	88.0	19.0	0
65	0.0	46.0	88.0	19.0	0
66	0.0	46.0	88.0	19.0	0
67	0.0	46.0	88.0	19.0	0
68	0.0	46.0	88.0	19.0	0
69	0.0	46.0	88.0	19.0	0
70	0.0	46.0	88.0	19.0	0
71	0.0	46.0	88.0	19.0	0
72	0.0	46.0	88.0	19.0	0
73	0.0	46.0	88.0	19.0	0
74	0.0	46.0	88.0	19.0	0
75	0.0	46.0	88.0	19.0	0
76	0.0	46.0	88.0	19.0	0
77	0.0	46.0	88.0	19.0	0
78	0.0	46.0	88.0	19.0	0
79	0.0	46.0	88.0	19.0	0
80	0.0	46.0	88.0	19.0	0
81	0.0	46.0	88.0	19.0	0
82	0.0	46.0	88.0	19.0	0
83	0.0	46.0	88.0	19.0	0
84	0.0	46.0	88.0	19.0	0
85	0.0	46.0	88.0	19.0	0
86	0.0	46.0	88.0	19.0	0
87	0.0	46.0	88.0	19.0	0
88	0.0	46.0	88.0	19.0	0
89	0.0	46.0	88.0	19.0	0
90	0.0	46.0	88.0	19.0	0
91	0.0	46.0	88.0	19.0	0
92	0.0	46.0	88.0	19.0	0
93	0.0	46.0	88.0	19.0	0
94	0.0	46.0	88.0	19.0	0
95	0.0	46.0	88.0	19.0	0
96	0.0	46.0	88.0	19.0	0
97	0.0	46.0	88.0	19.0	0
98	0.0	46.0	88.0	19.0	0
99	0.0	46.0	88.0	19.0	0
100	0.0	46.0	88.0	19.0	0
101	0.0	46.0	88.0	19.0	0
102	0.0	46.0	88.0	19.0	0
103	0.0	46.0	88.0	19.0	0
104	0.0	46.0	88.0	19.0	0
105	0.0	46.0	88.0	19.0	0
106	0.0	46.0	88.0	19.0	0
107	0.0	46.0	88.0	19.0	0
108	0.0	46.0	88.0	19.0	0
109	0.0	46.0	88.0	19.0	0
110	0.0	46.0	88.0	19.0	0
111	0.0	46.0	88.0	19.0	0
112	0.0	46.0	88.0	19.0	0
113	0.0	46.0	88.0	19.0	0
114	0.0	46.0	88.0	19.0	0
115	0.0	46.0	88.0	19.0	0
116	0.0	46.0	88.0	19.0	0
117	0.0	46.0	88.0	19.0	0
118	0.0	46.0	88.0	19.0	0
119	0.0	46.0	88.0	19.0	0
120	0.0	46.0	88.0	19.0	0
121	0.0	46.0	88.0	19.0	0
122	0.0	46.0	88.0	19.0	0
123	0.0	46.0	88.0	19.0	0
124	0.0	46.0	88.0	19.0	0
125	0.0	46.0	88.0	19.0	0
126	0.0	46.0	88.0	19.0	0
127	0.0	46.0	88.0	19.0	0
128	0.0	46.0	88.0	19.0	0
129	0.0	46.0	88.0	19.0	0
130	0.0	46.0	88.0	19.0	0
131	0.0	46.0	88.0	19.0	0
132	0.0	46.0	88.0	19.0	0
133	0.0	46.0	88.0	19.0	0
134	0.0	46.0	88.0	19.0	0
135	0.0	46.0	88.0	19.0	0
136	0.0	46.0	88.0	19.0	0
137	0.0	46.0	88.0	19.0	0
138	0.0	46.0	88.0	19.0	0
139	0.0	46.0	88.0	19.0	0
140	0.0	46.0	88.0	19.0	0
141	0.0	46.0	88.0	19.0	0
142	0.0	46.0	88.0	19.0	0
143	0.0	46.0	88.0	19.0	0
144	0.0	46.0	88.0	19.0	0
145	0.0	46.0	88.0	19.0	0
14					