# Directions for Unsupervised Machine Learning K-Means and DBSCAN



Image from: [[https://www.kaggle.com/nareshbhat/outlier-the-silent-killer (https://www.kaggle.com/nareshbhat/outlier-the-silent-killer)](https://www.kaggle.com/nareshbhat/outlier-the-silent-killer) (Links to an external site.)]

## Anomaly Detection Using Unsupervised Algorithms Overview

### Anomalies or outliers

Anomalies and outliers are often used interchangeably. Outliers are extreme data values that are distant from other values in your dataset. Outliers can distort the conclusions you draw from descriptive statistics (like mean and standard deviation) and from data visuals (like histograms and scatter plots). Many machine learning algorithms and statistical models are also sensitive to the range and distribution of data values used in your training data. Some of the most common causes of outliers in a dataset include data entry errors (human errors), measurement errors (instrument errors), and data processing errors (data manipulation errors). Consequently, outliers can come from multiple sources and hide in a single feature or n-dimensional feature space. Aside from the specific task of investigating these extreme data values (for anomaly detection) or in the case of naturally occurring outliers (novelties), the goal is often to remove or ignore outliers in your analysis.

**The following Anomaly Detection Techniques are good examples that can be used for Customer Segmentation.**

We all know that human diversity, whether it be gender, age, education, gender expression, or even food preferences, is a quality to be celebrated. While this is easier to accomplish in our dayto- day activities, from a business perspective, treating everyone as an individual is unrealistic, if not impossible, even though the business community is working on it. This is what motivates market segmentation. While you may not be able to serve your clients on an individual basis right now, you can do the next best thing by identifying their commonalities and differences, grouping them into small subgroups, and catering to their requirements. Even if ten thousand of your other clients received the same discount code in their email, this could give them the impression that you are communicating to them as an individual.

### Clustering

Clustering is an example of an unsupervised machine learning algorithm. When new data is presented, the fundamental idea behind employing a clustering algorithm for anomaly detection is to group data based on the structure and variables without any prior knowledge about the data. Another way to say this is that the purpose is to find "natural" groups in an unlabeled dataset. This information is then used to determine whether or not there are any anomalous points present. Customer segmentation, text clustering, and image segmentation are the most popular clustering techniques. There are many clustering algorithms. It is important to choose the best algorithm, and the best one to use depends on the database, client requirements, and client expectations.

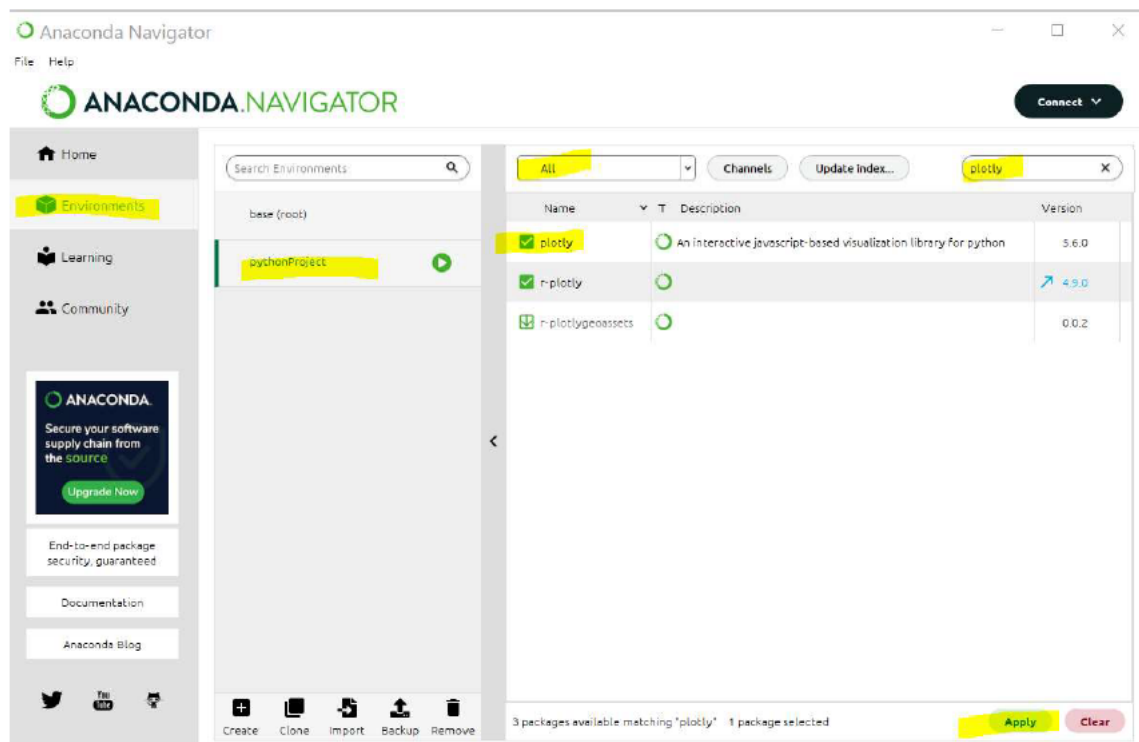### How do you segment your market?

K-means clustering is usually one of the first methods used. Because k-means is one of, if not the only, data clustering techniques taught in statistics classrooms, it's only natural that business analysts would use it when a project requires market segmentation. As we learned in the lecture on k-means, k-means is simple to learn and effective with massive data sets. But, like all statistical approaches, k-means clustering relies on specific assumptions, and if those assumptions prove correct, the method is useless. Every data point is given to a cluster in the k-means clustering result, and the mean of the points is the centroid of each cluster. As a result, each point contributed to the cluster's formation. In reality, it's not uncommon for certain data points to stray from the "norm", resulting in anomalies or outliers. A single outlier in a dataset can cause the entire clustering result to be skewed. Furthermore, clusters formed using k-means are more spherical, yet, inaccuracies can develop if clusters are of other odd or elliptical shapes. Scree plots and silhouette plots can help you decide on parameter k; however, model sensitivity to anomalies significantly influences them. The goal of developing a new strategy is to overcome the drawbacks of existing methods. The majority of k-means drawbacks are dealt with or handled by DBSCAN. DBSCAN is a density-based non-parametric unsupervised learning algorithm. The DBSCAN has fewer assumptions, a more flexible model, and a larger range of applications than k-means. However, DBSCAN, like k-means, groups data according to their similarities using distance functions. DBSCAN also uses densities to form groups of data. The term "density" refers to the number of points in the region indicated by the model parameters.

## Let's Get Started with Unsupervised Machine Learning

In the assignment, we will focus k-means and DBSCAN algorithms by using the Scikit-Learn Python library. From this dataset, we need to calculate some patterns; as it is an unsupervised method, we don't know what to calculate exactly. The column names will help you make sense of the data.

## STEP 0: Install plotly by going to Anaconda Navigator and selecting the

Environments Tab. Select pythonProject, All packages, and search for plotly. Select plotly, and select Apply to start the installation. Once it has been installed go to Home and open Jupytner Notebook. Next, create a new Jupyter Notebook file, Name the file, and then Enter your name and date at the top of the file. SAVE file.



## STEP 1: Import Libraries

```
In [163]:    1  # Import Python Libraries
             2  import pandas as pd
             3  import seaborn as sns
             4  import matplotlib.pyplot as plt
             5  import numpy as np
             6  from scipy import stats
             7  from sklearn.metrics import silhouette_score
             8  from sklearn.cluster import DBSCAN
             9  from sklearn.cluster import KMeans
            10  from itertools import product
            11  from mpl_toolkits.mplot3d import Axes3D
            12  import plotly as py
            13  import plotly.graph_objs as go
            14  #filter warnings
            15  import warnings
            16  warnings.filterwarnings("ignore")
            17
            18  plt.rcParams.update({'font.size': 8})
```

## STEP 2: Load Dataset & Get Dataset Shape

```
In [164]:    1  mall_data = pd.read_csv('Mall_Customers.csv')
             2  print('There are {} rows and {} columns in our dataset'.format(mall_data.shape[0],mall_data.shape[1]))
             3
             4  mall_data.head()
             5  mall_data.describe()
             6
```

There are 200 rows and 5 columns in our dataset

Out[164]:

|  | CustomerID | Age | Annual Income (k$) | Spending Score (1-100) |
|---|---|---|---|---|
| count | 200.000000 | 200.000000 | 200.000000 | 200.000000 |
| mean | 100.500000 | 38.850000 | 60.560000 | 50.200000 |
| std | 57.879185 | 13.969007 | 26.264721 | 25.823522 |
| min | 1.000000 | 18.000000 | 15.000000 | 1.000000 |
| 25% | 50.750000 | 28.750000 | 41.500000 | 34.750000 |
| 50% | 100.500000 | 36.000000 | 61.500000 | 50.000000 |
| 75% | 150.250000 | 49.000000 | 78.000000 | 73.000000 |
| max | 200.000000 | 70.000000 | 137.000000 | 99.000000 |

## STEP 3: Clean the data

```
In [165]:   1  # We see there are no missing data points
            2  mall_data.isnull().sum()
```

```
Out[165]:  CustomerID              0
           Genre                   0
           Age                     0
           Annual Income (k$)      0
           Spending Score (1-100)  0
           dtype: int64
```
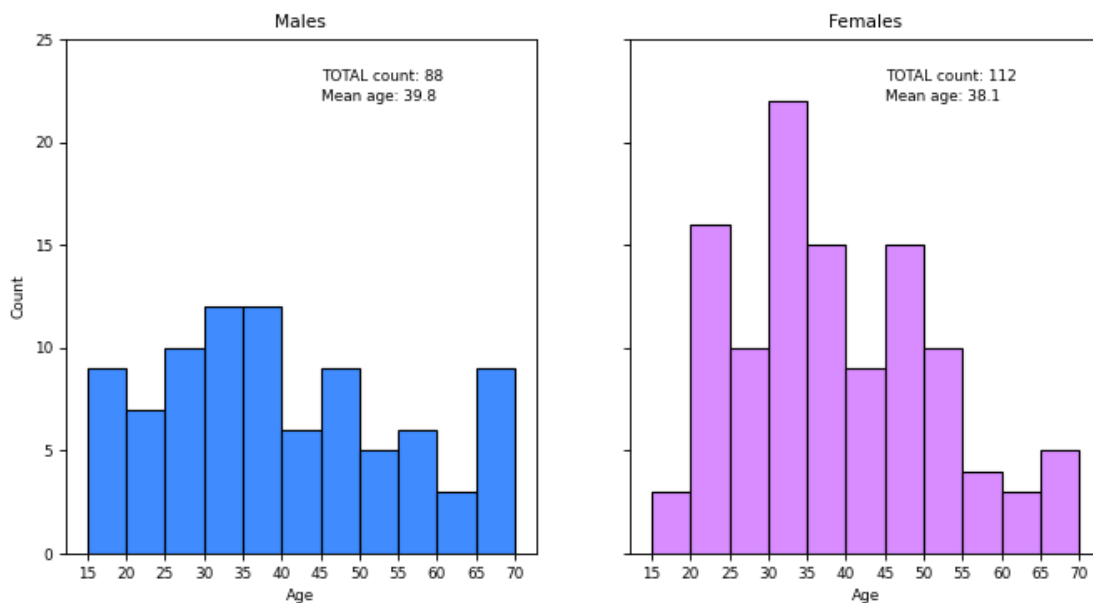
## STEP 4: Performing the Exploratory Data Analysis (EDA)

### Distribution - Age

```
In [166]:   1  males_age = mall_data[mall_data['Genre']=='Male']['Age'] # subset with males age
            2  females_age = mall_data[mall_data['Genre']=='Female']['Age'] # subset with females age
            3  age_bins = range(15,75,5)
            4  # males histogram
            5  fig2, (ax1, ax2) = plt.subplots(1, 2, figsize=(10,5), sharey=True)
            6  sns.histplot(males_age,bins=age_bins, color='#0066ff', ax=ax1)
            7  ax1.set_xticks(age_bins)
            8  ax1.set_ylim(top=25)
            9  ax1.set_title('Males')
           10  ax1.set_ylabel('Count')
           11  ax1.text(45,23, "TOTAL count: {}".format(males_age.count()))
           12  ax1.text(45,22, "Mean age: {:.1f}".format(males_age.mean()))
           13  # females histogram
           14  sns.histplot(females_age, bins=age_bins, color='#cc66ff', ax=ax2, )
           15  ax2.set_xticks(age_bins)
           16  ax2.set_title('Females')
           17  ax2.set_ylabel('Count')
           18  ax2.text(45,23, "TOTAL count: {}".format(females_age.count()))
           19  ax2.text(45,22, "Mean age: {:.1f}".format(females_age.mean()))
           20  plt.show()
           21  print('Kolgomorov-Smirnov test p-value:{:.2f}'.format(stats.ks_2samp(males_age, females_age)[1]))
```

Figure 103



```
Kolgomorov-Smirnov test p-value:0.49
```

Male age is more evenly distributed than females. The biggest age group for women is between 30-35, and for men, it is between 30-40. The Kolmogorov-Smirnov test shows the difference between the males and females age groups is statistically insignificant as it calculates the P-value. A P-value that is less than or equal to 0.05 is statistically significant.
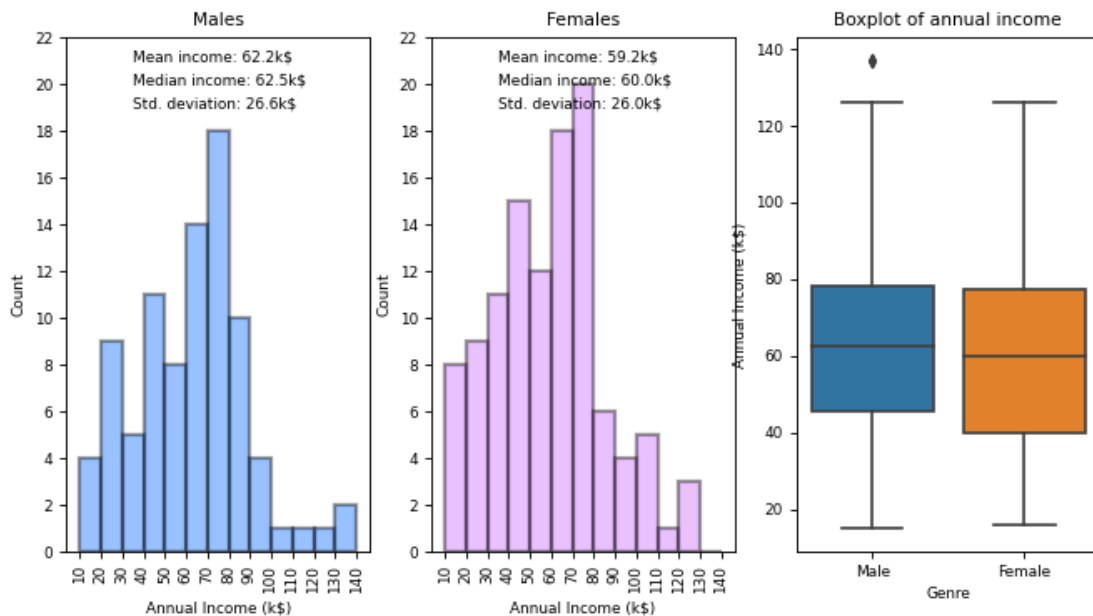
### Distribution - Annual Income (k$)

```
In [167]:     1  males_income = mall_data[mall_data['Genre']=='Male']['Annual Income (k$)']
              2  # subset with males income
              3  females_income = mall_data[mall_data['Genre']=='Female']['Annual Income (k$)']
              4  # subset with females income
              5
              6  my_bins = range(10,150,10)
              7  # males histogram
              8  fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(10,5))
              9  sns.distplot(males_income, bins=my_bins, kde=False, color='#0066ff', ax=ax1,
             10  hist_kws=dict(edgecolor="k", linewidth=2))
             11
             12  ax1.set_xticks(my_bins)
             13  ax1.set_xticklabels(ax1.get_xticks(), rotation = 90)
             14
             15  ax1.set_yticks(range(0,24,2))
             16  ax1.set_ylim(0,22)
             17  ax1.set_title('Males')
             18  ax1.set_ylabel('Count')
             19  ax1.text(35,21, "Mean income: {:.1f}k$".format(males_income.mean()))
             20  ax1.text(35,20, "Median income: {:.1f}k$".format(males_income.median()))
             21  ax1.text(35,19, "Std. deviation: {:.1f}k$".format(males_income.std()))
             22  # females histogram
             23  sns.distplot(females_income, bins=my_bins, kde=False, color='#cc66ff',
             24  ax=ax2, hist_kws=dict(edgecolor="k", linewidth=2))
             25  ax2.set_xticks(my_bins)
             26  ax2.set_xticklabels(ax2.get_xticks(), rotation = 90)
             27
             28  ax2.set_yticks(range(0,24,2))
             29  ax2.set_ylim(0,22)
             30  ax2.set_title('Females')
             31  ax2.set_ylabel('Count')
             32  ax2.text(35,21, "Mean income: {:.1f}k$".format(females_income.mean()))
             33  ax2.text(35,20, "Median income: {:.1f}k$".format(females_income.median()))
             34  ax2.text(35,19, "Std. deviation: {:.1f}k$".format(females_income.std()))
             35  # boxplot
             36  sns.boxplot(x='Genre', y='Annual Income (k$)', data=mall_data, ax=ax3)
             37  ax3.set_title('Boxplot of annual income')
             38  plt.show()
             39  print('Kolgomorov-Smirnov test p-value:{:.2f}'.format(stats.ks_2samp(males_income, females_income)[1]))
```

Figure 104
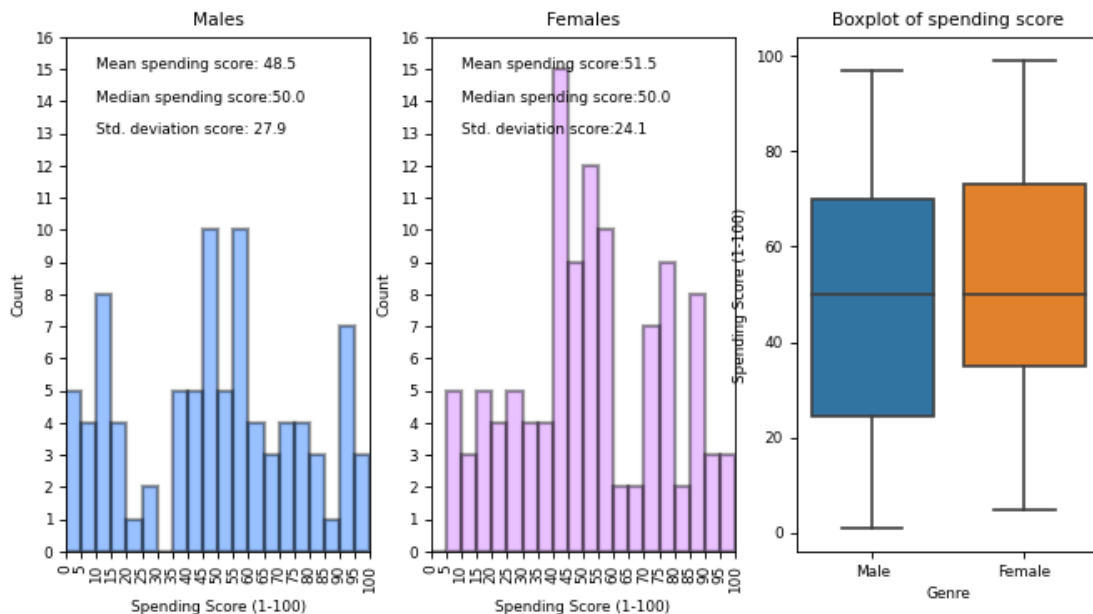


Kolgomorov-Smirnov test p-value:0.78

**Distribution - Spending Score (1-100)**

```python
1  males_spending = mall_data[mall_data['Genre']=='Male']['Spending Score (1-100)']
2  # subset with males age
3  females_spending = mall_data[mall_data['Genre']=='Female']['Spending Score (1-100)']
4  # subset with females age
5  spending_bins = range(0,105,5)
6  # males histogram
7  fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(10,5))
8  sns.distplot(males_spending, bins=spending_bins, kde=False, color='#0066ff',
9  ax=ax1, hist_kws=dict(edgecolor="k", linewidth=2))
10 ax1.set_xticks(spending_bins)
11 ax1.set_xticklabels(ax1.get_xticks(), rotation = 90)
12
13 ax1.set_xlim(0,100)
14 ax1.set_yticks(range(0,17,1))
15 ax1.set_ylim(0,16)
16 ax1.set_title('Males')
17 ax1.set_ylabel('Count')
18 ax1.text(10,15, "Mean spending score: {:.1f}".format(males_spending.mean()))
19 ax1.text(10,14, "Median spending score:{:.1f}".format(males_spending.median()))
20 ax1.text(10,13, "Std. deviation score: {:.1f}".format(males_spending.std()))
21 # females histogram
22 sns.distplot(females_spending, bins=spending_bins, kde=False,
23 color='#cc66ff', ax=ax2, hist_kws=dict(edgecolor="k", linewidth=2))
24 ax2.set_xticks(spending_bins)
25 ax2.set_xticklabels(ax2.get_xticks(), rotation = 90)
26
27 ax2.set_xlim(0,100)
28 ax2.set_yticks(range(0,17,1))
29 ax2.set_ylim(0,16)
30 ax2.set_title('Females')
31 ax2.set_ylabel('Count')
32 ax2.text(10,15, "Mean spending score:{:.1f}".format(females_spending.mean()))
33 ax2.text(10,14, "Median spending score:{:.1f}".format(females_spending.median()))
34 ax2.text(10,13, "Std. deviation score:{:.1f}".format(females_spending.std()))
35 # boxplot
36 sns.boxplot(x='Genre', y='Spending Score (1-100)', data=mall_data, ax=ax3)
37 ax3.set_title('Boxplot of spending score')
38 plt.show()
39 print('Kolgomorov-Smirnov test p-value:{:.2f}'.format(stats.ks_2samp(males_spending, females_spending)[1]))
```
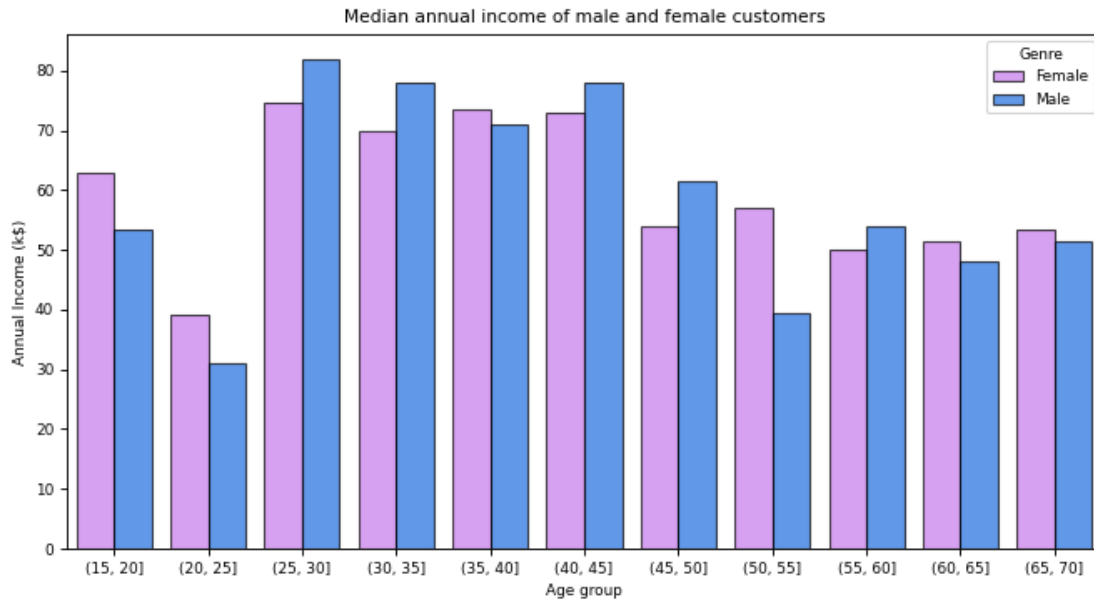
Figure 105



Kolgomorov-Smirnov test p-value:0.29

**Distribution - Median Annual Income of Male and Female Customers**

```
1  medians_by_age_group = mall_data.groupby(["Genre",pd.cut(mall_data['Age'], age_bins)]).median()
2  medians_by_age_group.index = medians_by_age_group.index.set_names(['Genre', 'Age_group'])
3  medians_by_age_group.reset_index(inplace=True)
4  fig, ax = plt.subplots(figsize=(10,5))
5  sns.barplot(x='Age_group', y='Annual Income (k$)', hue='Genre',
6              data=medians_by_age_group, palette=['#cc66ff','#0066ff'],
7              alpha=0.7,edgecolor='k',ax=ax)
8  ax.set_title('Median annual income of male and female customers')
9  ax.set_xlabel('Age group')
10 plt.show()
```
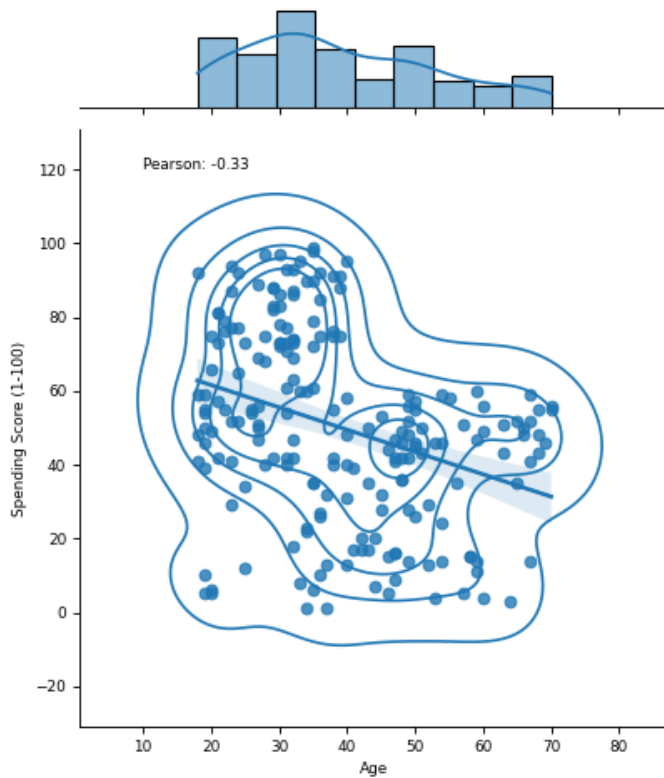
Figure 106



**Correlations**

**Pearson's Correlation for Age & Spending Score (1-100) with jointplot**

```
1  medians_by_age_group = mall_data.groupby(["Genre",pd.cut(mall_data['Age'], age_bins)]).median()
2  medians_by_age_group.index = medians_by_age_group.index.set_names(['Genre', 'Age_group'])
3  medians_by_age_group.reset_index(inplace=True)
4  fig, ax = plt.subplots(figsize=(10,5))
5  sns.barplot(x='Age_group', y='Annual Income (k$)', hue='Genre',
6              data=medians_by_age_group, palette=['#cc66ff','#0066ff'],
7              alpha=0.7,edgecolor='k',ax=ax)
8  ax.set_title('Median annual income of male and female customers')
9  ax.set_xlabel('Age group')
10 plt.show()
```

```
In [197]:  1  from scipy.stats import pearsonr
           2  # calculating Pearson's correlation
           3  corr, _ = pearsonr(mall_data['Age'], mall_data['Spending Score (1-100)'])
           4  jp = sns.jointplot(x='Age', y='Spending Score (1-100)', data=mall_data, kind='reg').plot_joint(
           5                      sns.kdeplot, zorder=0, n_levels=6)
           6
           7  plt.text(10,120, 'Pearson: {:.2f}'.format(corr))
           8  plt.show()
```
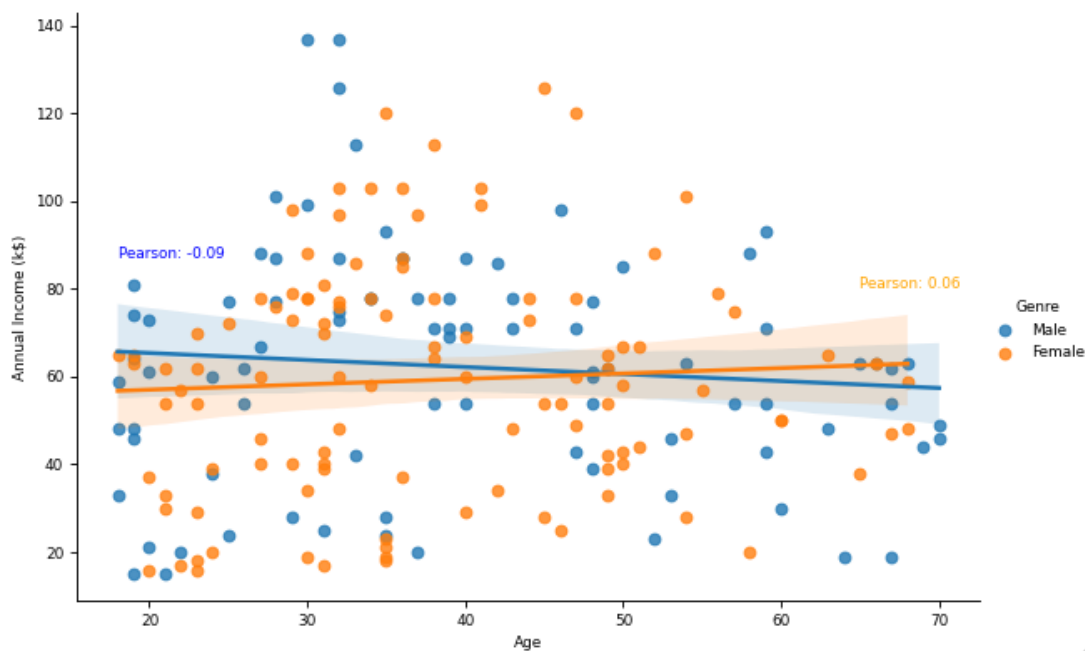
Figure 120



**Pearson's Correlation for Age & Annual Income with lineplot (lmplot)**

```
In [202]:  1  # calculating Pearson's correlations
           2  corr1, _ = pearsonr(males_age.values, males_income.values)
           3  corr2, _ = pearsonr(females_age.values, females_income.values)
           4  sns.lmplot(x='Age', y='Annual Income (k$)', data=mall_data, hue='Genre',aspect=1.5)
           5  plt.text(18,87, 'Pearson: {:.2f}'.format(corr1), color='blue')
           6  plt.text(65,80, 'Pearson: {:.2f}'.format(corr2), color='orange')
           7  plt.show()
```
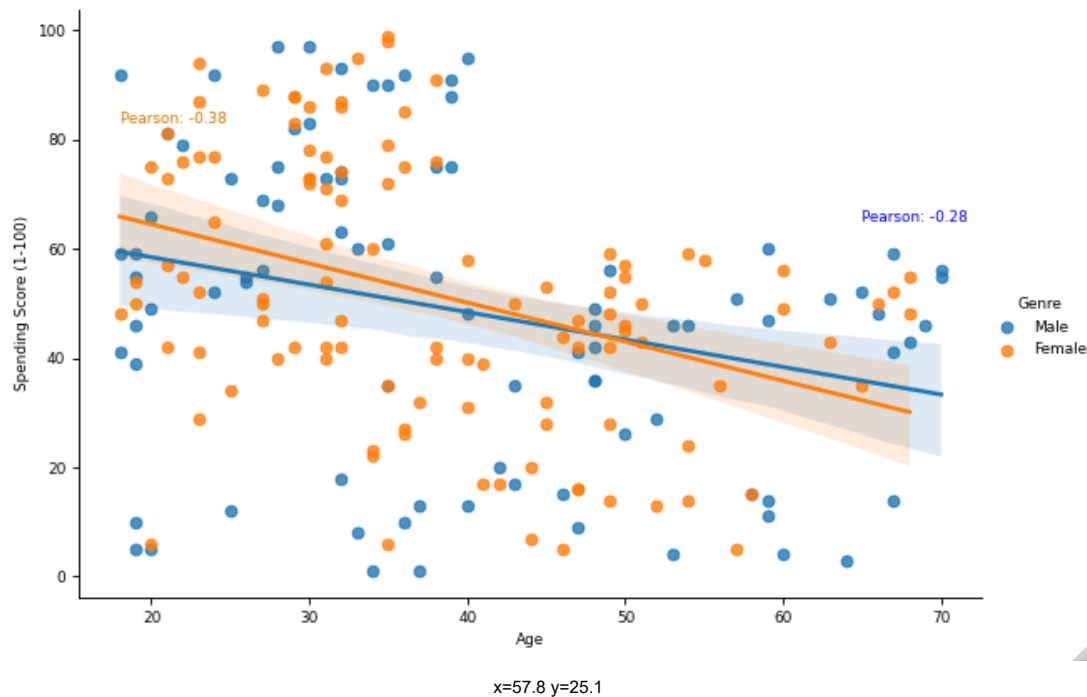
Figure 125

**Pearson's Correlation for Age & Spending Score (1-100) with line plot (lmplot)**

```
In [203]:    1  # calculating Pearson's correlations
             2  corr1, _ = pearsonr(males_age.values, males_spending.values)
             3  corr2, _ = pearsonr(females_age.values, females_spending.values)
             4  sns.lmplot(x='Age', y='Spending Score (1-100)', data=mall_data, hue='Genre', aspect=1.5)
             5  plt.text(65,65, 'Pearson: {:.2f}'.format(corr1), color='blue')
             6  plt.text(18,83, 'Pearson: {:.2f}'.format(corr2), color='#d97900')
             7  plt.show()
```
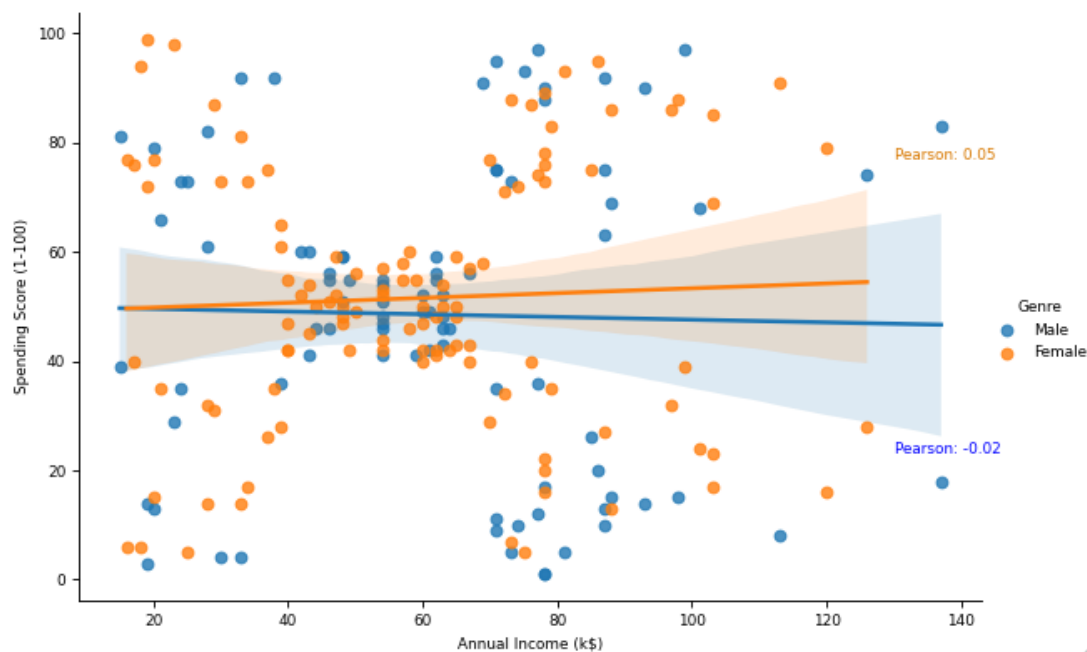
Figure 126



x=57.8 y=25.1

**Pearson's Correlation for Annual Income & Spending Score (1-100) with line plot (lmplot)**

```
In [173]:    1  # calculating Pearson's correlations
             2  corr1, _ = pearsonr(males_income.values, males_spending.values)
             3  corr2, _ = pearsonr(females_income.values, females_spending.values)
             4  sns.lmplot(x='Annual Income (k$)', y='Spending Score (1-100)', data=mall_data,
             5  hue='Genre', aspect=1.5)
             6  plt.text(130,23, 'Pearson: {:.2f}'.format(corr1), color='blue')
             7  plt.text(130,77, 'Pearson: {:.2f}'.format(corr2), color='#d97900')
             8  plt.show()
```

Figure 110

## STEP 5: K-Means

The k-means algorithm is the most used partitional clustering technique. It gained popularity since it is so easy and simple to implement, and it has had many successes in many diverse sectors. As we learned in our lecture for this module, k-means is a greedy algorithm. Remember with the greedy algorithm to be aware of global optimization. Also, remember that the Euclidean distance is implemented, and the number of clusters has to be predefined. In this assignment, you will use the Elbow Method and Silhouette Score to calculate the k-value. For clustering, you will only use numerical columns.

```
In [174]:   1  X_numerics = mall_data[['Age', 'Annual Income (k$)', 'Spending Score (1-100)']]
            2  # subset with numeric variables only
```
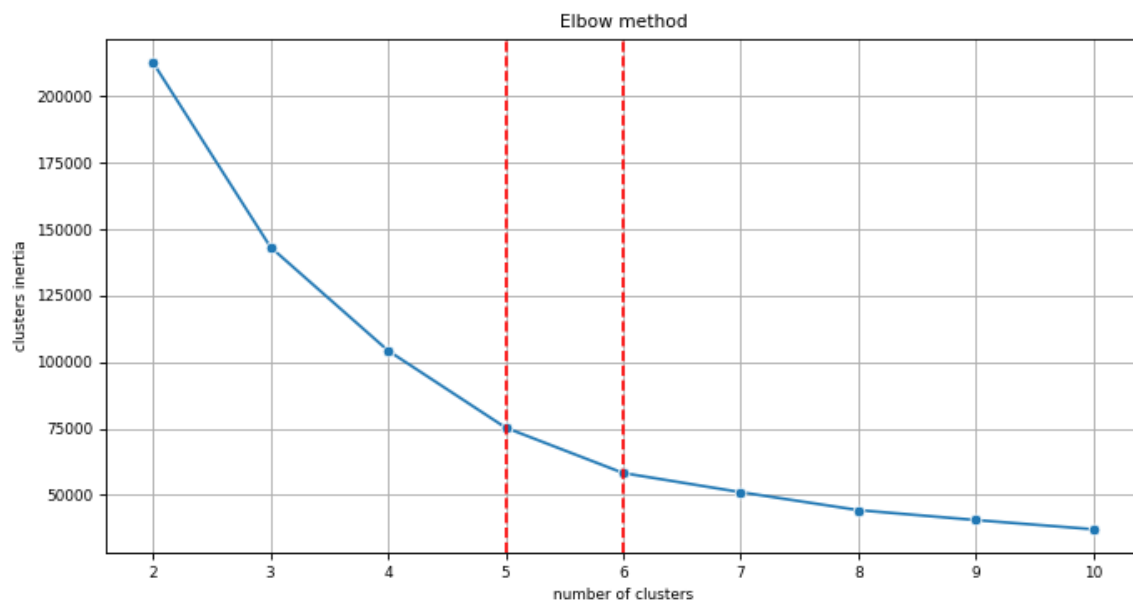
**Elbow Method**

In order to find the optimal K number, look at the "elbow" in the graph.

```
In [175]:   1  n_clusters = [2,3,4,5,6,7,8,9,10] # number of clusters
            2  clusters_inertia = [] # inertia of clusters
            3  s_scores = [] # silhouette scores
            4  for n in n_clusters:
            5      KM_est = KMeans(n_clusters=n, init='k-means++').fit(X_numerics)
            6      clusters_inertia.append(KM_est.inertia_) # data for the elbow method
            7      silhouette_avg = silhouette_score(X_numerics, KM_est.labels_)
            8      s_scores.append(silhouette_avg) # data for the silhouette score method
```

```
In [176]:   1  fig, ax = plt.subplots(figsize=(10,5))
            2  ax = sns.lineplot(x=n_clusters, y=clusters_inertia, marker='o', ax=ax)
            3  ax.set_title("Elbow method")
            4  ax.set_xlabel("number of clusters")
            5  ax.set_ylabel("clusters inertia")
            6  ax.axvline(5, ls="--", c="red")
            7  ax.axvline(6, ls="--", c="red")
            8  plt.grid()
            9  plt.show()
```
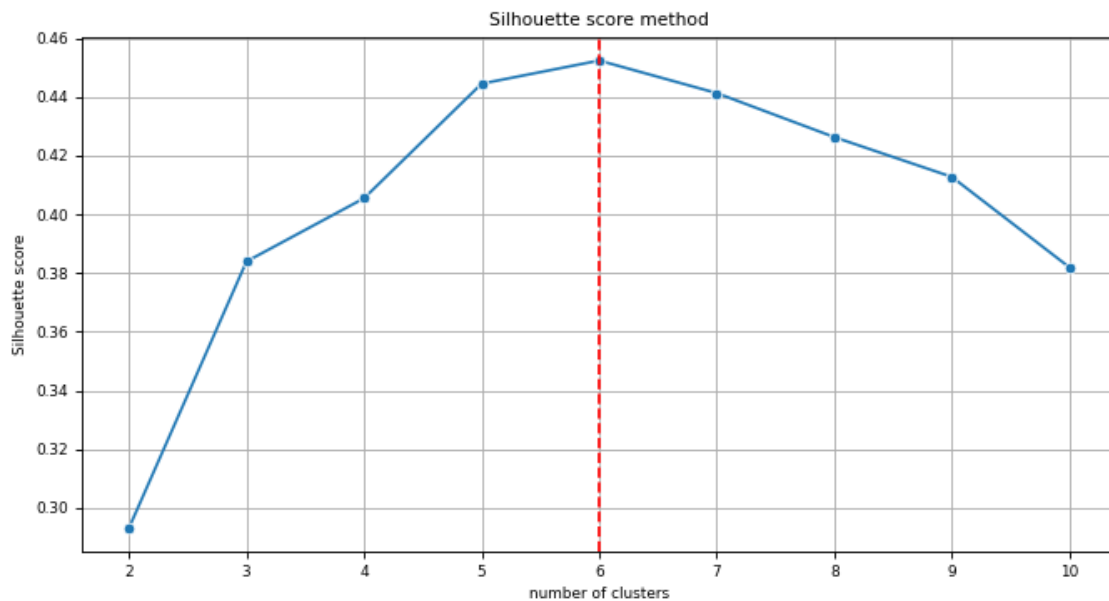
Figure 111



**Silhouette Score Method**

Since there is no clear "elbow", the Silhouette Score will also be utilized. The Silhouette Score illustrates the distance between clusters, ranging between -1 and 1, and the best score is the number closest to 1.

```
1  fig, ax = plt.subplots(figsize=(10,5))
2  ax = sns.lineplot(x=n_clusters, y=s_scores, marker='o', ax=ax)
3  ax.set_title("Silhouette score method")
4  ax.set_xlabel("number of clusters")
5  ax.set_ylabel("Silhouette score")
6  ax.axvline(6, ls="--", c="red")
7  plt.grid()
8  plt.show()
```
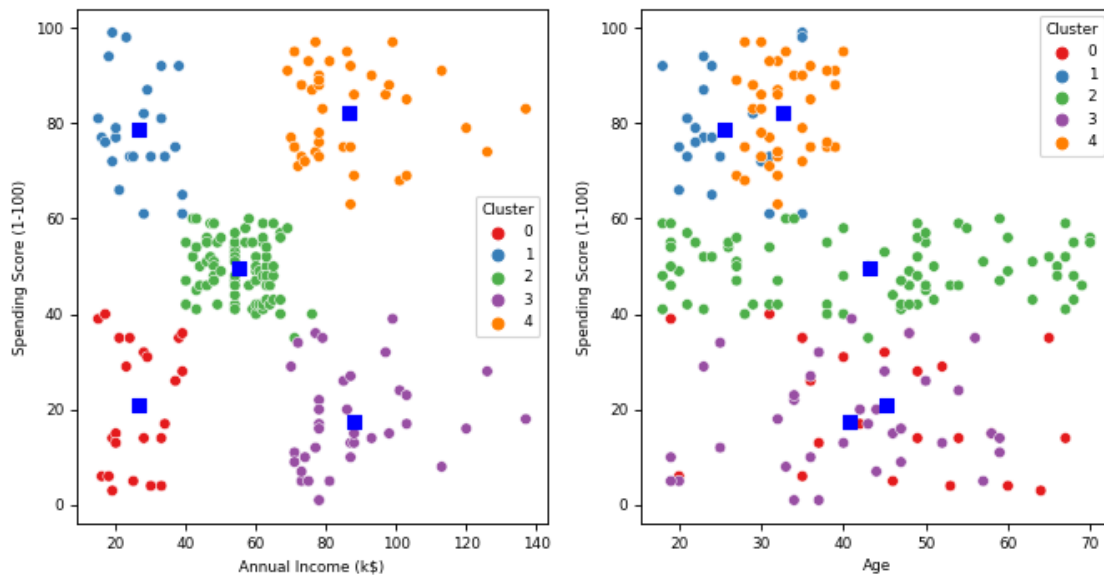
Figure 112



**Let's try 5 Clusters**

```
1   # initialize and fit K-Means model
2   KM_5_clusters = KMeans(n_clusters=5, init='k-means++').fit(X_numerics)
3   KM5_clustered = X_numerics.copy()
4   # append labels to points
5   KM5_clustered.loc[:,'Cluster'] = KM_5_clusters.labels_
6   fig1, (axes) = plt.subplots(1,2,figsize=(10,5))
7   scat_1 = sns.scatterplot(x='Annual Income (k$)', y='Spending Score (1-100)',
8   data=KM5_clustered,hue='Cluster', ax=axes[0], palette='Set1', legend='full')
9   sns.scatterplot(x='Age', y='Spending Score (1-100)', data=KM5_clustered,
10  hue='Cluster', palette='Set1', ax=axes[1], legend='full')
11  axes[0].scatter(KM_5_clusters.cluster_centers_[:,1],KM_5_clusters.cluster_centers_[:,2], marker='s', s=40, c="blue")
12  axes[1].scatter(KM_5_clusters.cluster_centers_[:,0],KM_5_clusters.cluster_centers_[:,2], marker='s', s=40, c="blue")
13  plt.show()
```

Figure 113

```
In [179]:    1  # Size of clusters
             2  KM_clust_sizes = KM5_clustered.groupby('Cluster').size().to_frame()
             3  KM_clust_sizes.columns = ["KM_size"]
             4  KM_clust_sizes
```
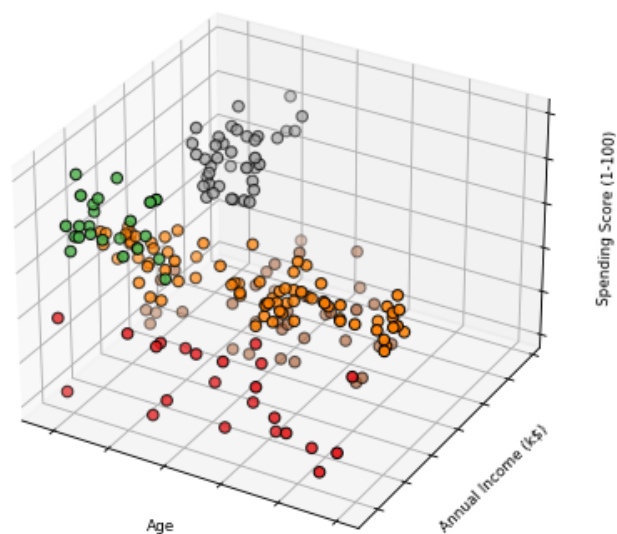
Out[179]:

|         | KM_size |
|---------|---------|
| Cluster |         |
| 0       | 23      |
| 1       | 23      |
| 2       | 79      |
| 3       | 36      |
| 4       | 39      |

```
In [180]:    1  # if not installed
             2  #!pip install ipympl
             3
             4  #Enable Javascript in your browser
             5
             6  #Create a 3D projection of 5 generated clusters
             7  from mpl_toolkits.mplot3d import Axes3D
             8  # for creating a responsive plot
             9  %matplotlib widget
            10
            11  # creating figure
            12  fig = plt.figure(figsize=(7, 7))
            13  ax=fig.add_subplot(111,projection='3d')
            14
            15  ax.scatter(KM5_clustered['Age'],
            16      KM5_clustered['Annual Income (k$)'],
            17      KM5_clustered['Spending Score (1-100)'],
            18      c=KM5_clustered['Cluster'],
            19      s=35, edgecolor='k', cmap=plt.cm.Set1)
            20  ax.w_xaxis.set_ticklabels([])
            21  ax.w_yaxis.set_ticklabels([])
            22  ax.w_zaxis.set_ticklabels([])
            23  ax.set_xlabel('Age')
            24  ax.set_ylabel('Annual Income (k$)')
            25  ax.set_zlabel('Spending Score (1-100)')
            26  ax.set_title('3D view of K-Means 5 clusters')
            27  ax.dist = 12
            28  plt.show()
```
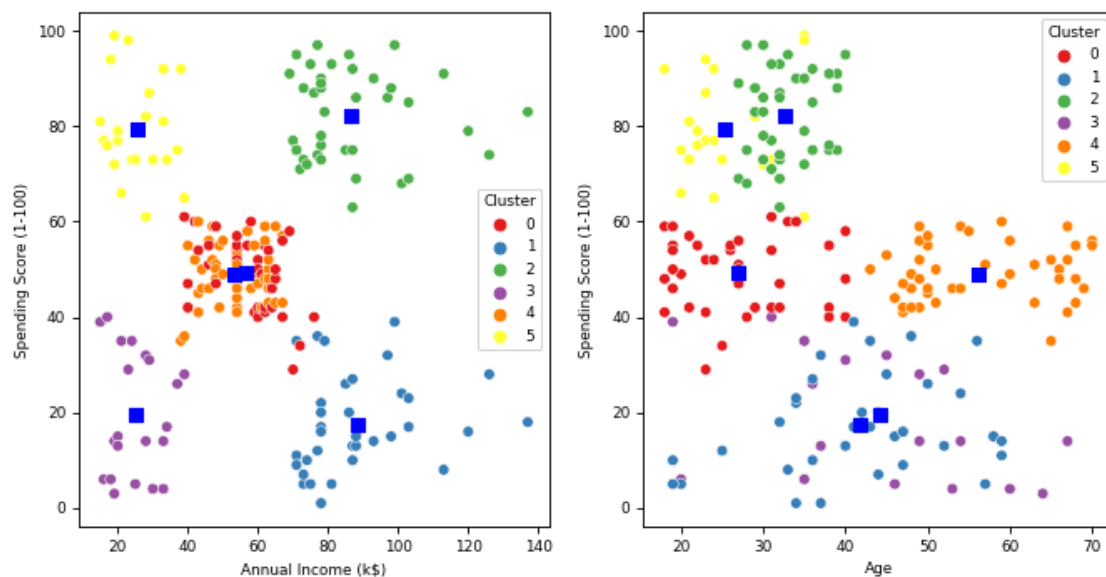
Figure 114



3D view of K-Means 5 clusters

**Let's try 6 Clusters**

In [181]:
```python
#initialize and fit K-Means model
KM_6_clusters = KMeans(n_clusters=6, init='k-means++').fit(X_numerics)
KM6_clustered = X_numerics.copy()
# append labels to points
KM6_clustered.loc[:,'Cluster'] = KM_6_clusters.labels_
fig11, (axes) = plt.subplots(1,2,figsize=(10,5))
sns.scatterplot(x='Annual Income (k$)', y='Spending Score (1-100)',
data=KM6_clustered,hue='Cluster', ax=axes[0], palette='Set1', legend='full')
sns.scatterplot(x='Age', y='Spending Score (1-100)', data=KM6_clustered,
hue='Cluster', palette='Set1', ax=axes[1], legend='full')
# plotting centroids
axes[0].scatter(KM_6_clusters.cluster_centers_[:,1],
KM_6_clusters.cluster_centers_[:,2], marker='s', s=40, c="blue")
axes[1].scatter(KM_6_clusters.cluster_centers_[:,0],
KM_6_clusters.cluster_centers_[:,2], marker='s', s=40, c="blue")
plt.show()
```

Figure 115



In [182]:
```python
# Size of Clusters
KM6_clust_sizes = KM6_clustered.groupby('Cluster').size().to_frame()
KM6_clust_sizes.columns = ["KM_size"]
KM6_clust_sizes
```
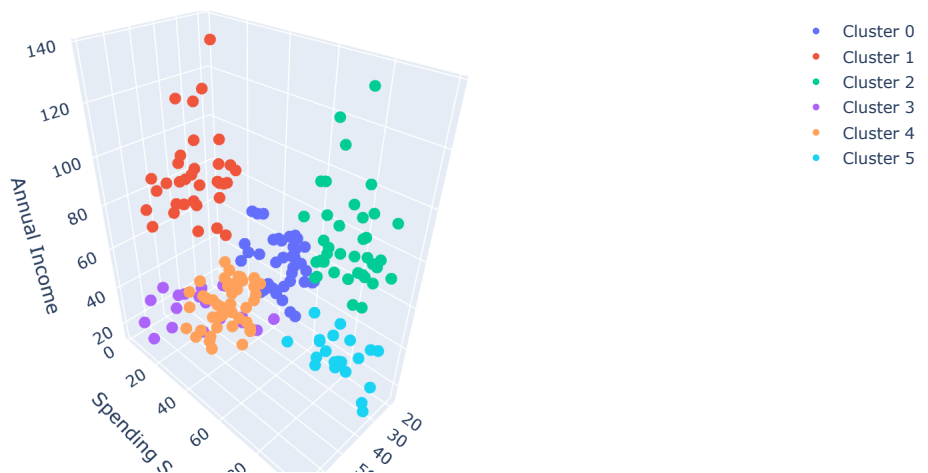
Out[182]:

| Cluster | KM_size |
|---------|---------|
| 0 | 38 |
| 1 | 35 |
| 2 | 39 |
| 3 | 21 |
| 4 | 45 |
| 5 | 22 |

```
In [204]:  1  # Create an interactive 3D view of K-Means with 6 clusters with Plotly
           2  def tracer(db, n, name):
           3      '''
           4      This function returns trace object for Plotly
           5      '''
           6      return go.Scatter3d(
           7          x = db[db['Cluster']==n]['Age'],
           8          y = db[db['Cluster']==n]['Spending Score (1-100)'],
           9          z = db[db['Cluster']==n]['Annual Income (k$)'],
          10          mode = 'markers',
          11          name = name,
          12          marker = dict(size = 5)
          13      )
          14  trace0 = tracer(KM6_clustered, 0, 'Cluster 0')
          15  trace1 = tracer(KM6_clustered, 1, 'Cluster 1')
          16  trace2 = tracer(KM6_clustered, 2, 'Cluster 2')
          17  trace3 = tracer(KM6_clustered, 3, 'Cluster 3')
          18  trace4 = tracer(KM6_clustered, 4, 'Cluster 4')
          19  trace5 = tracer(KM6_clustered, 5, 'Cluster 5')
          20  data = [trace0, trace1, trace2, trace3, trace4, trace5]
          21  layout = go.Layout(
          22      title = 'Clusters by K-Means',
          23      scene = dict(
          24          xaxis = dict(title = 'Age'),
          25          yaxis = dict(title = 'Spending Score'),
          26          zaxis = dict(title = 'Annual Income')
          27      )
          28  )
          29  fig = go.Figure(data=data, layout=layout)
          30  py.offline.iplot(fig)
```



Clusters by K-Means

## STEP 6: DBSCAN (Density-Based Spatial Clustering of Applications with

Noise) The basis of this algorithm is based on the concept of density. We will use the distance (eps) and the minimum number of points within the distance as parameters. Again, we will use the Euclidean distance for the DBSCAN. Create a matrix of investigated combinations and the number of generated clusters.

```
In [184]:  1  # create a matrix of investigated combinations.
           2  eps_values = np.arange(8,12.75,0.25) # eps values to be investigated
           3  min_samples = np.arange(3,10) # min_samples values to be investigated
           4  DBSCAN_params = list(product(eps_values, min_samples))
```

```
In [185]:  1  # collecting number of generated clusters
           2  no_of_clusters = []
           3  sil_score = []
           4  for p in DBSCAN_params:
           5      DBS_clustering = DBSCAN(eps=p[0], min_samples=p[1]).fit(X_numerics)
           6      no_of_clusters.append(len(np.unique(DBS_clustering.labels_)))
           7      sil_score.append(silhouette_score(X_numerics, DBS_clustering.labels_))
```
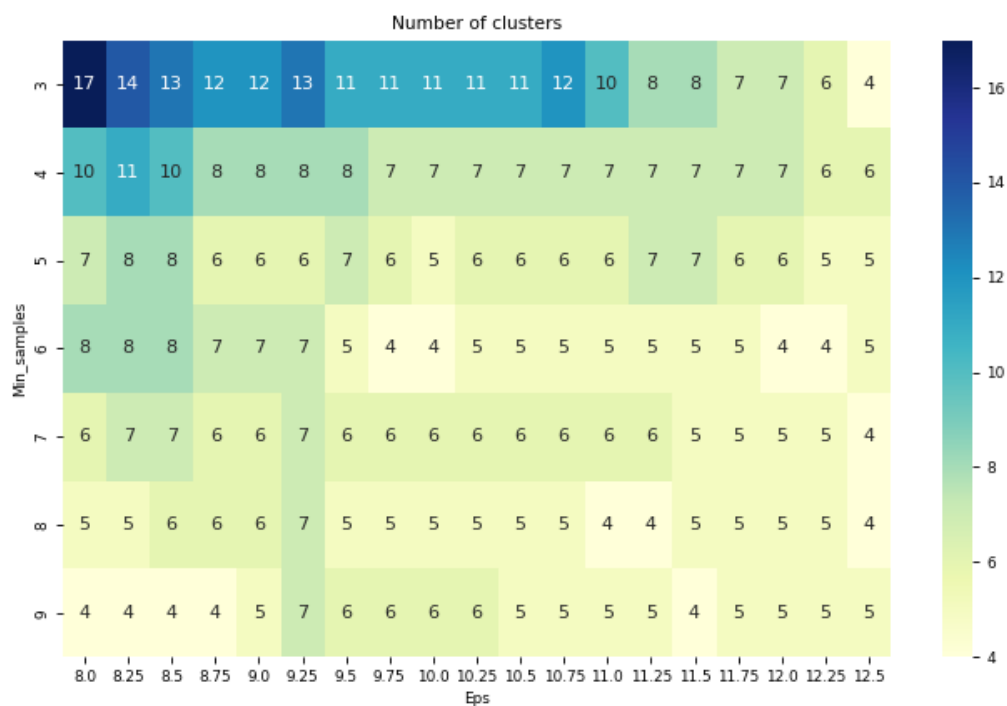
### Create Heatmaps

The heatplot illustrates how many clusters were generated by the algorithm with the respective parameters combinations created above. You can see there are clusters from 17 to 4.

```
In [186]:    1  tmp = pd.DataFrame.from_records(DBSCAN_params, columns =['Eps', 'Min_samples'])
             2  tmp['No_of_clusters'] = no_of_clusters
             3  pivot_1 = pd.pivot_table(tmp, values='No_of_clusters', index='Min_samples',
             4  columns='Eps')
             5  fig, ax = plt.subplots(figsize=(10,6))
             6  sns.heatmap(pivot_1, annot=True,annot_kws={"size": 10}, cmap="YlGnBu", ax=ax)
             7  ax.set_title('Number of clusters')
             8  plt.show()
```
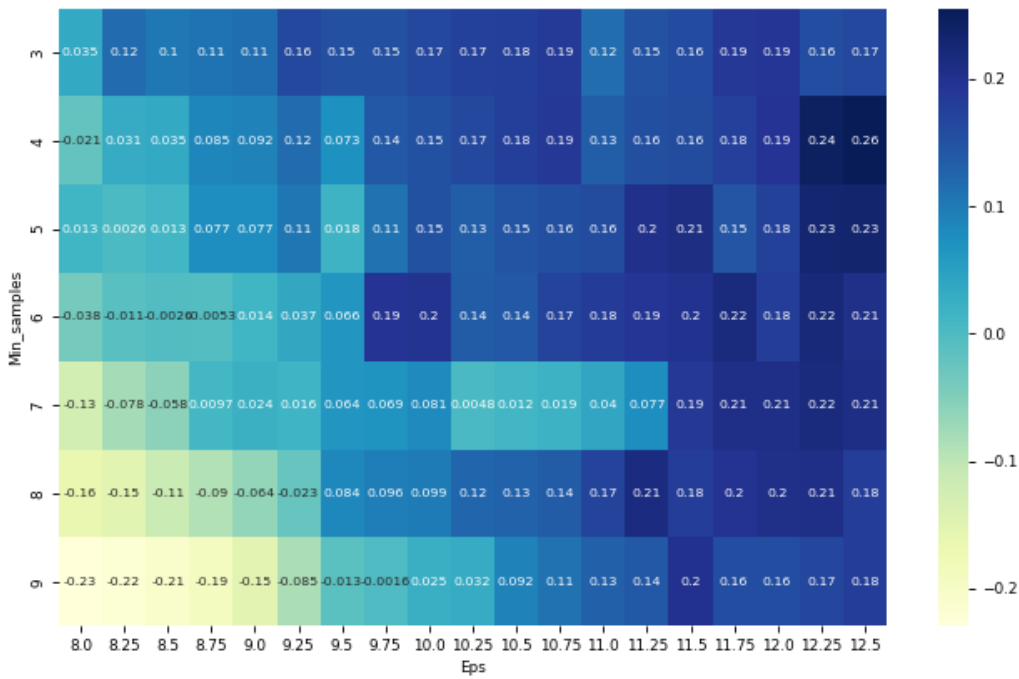
Figure 116



Number of clusters

```
In [187]:   1  # In the heatmap below, we see the Global maximum is 0.26
            2  # for eps=12.5 and
            3  min_samples=4.
            4  tmp = pd.DataFrame.from_records(DBSCAN_params, columns =['Eps', 'Min_samples'])
            5  tmp['Sil_score'] = sil_score
            6  pivot_1 = pd.pivot_table(tmp, values='Sil_score', index='Min_samples',
            7  columns='Eps')
            8  fig, ax = plt.subplots(figsize=(10,6))
            9  sns.heatmap(pivot_1, annot=True, annot_kws={"size": 7}, cmap="YlGnBu", ax=ax)
           10  plt.show()
```

Figure 117



```
In [188]:   1  # Based on the graph above let's use the Global
            2  # maximum settings (eps = 12.5 and min_samples = 4)
            3  DBS_clustering = DBSCAN(eps=12.5, min_samples=4).fit(X_numerics)
            4  DBSCAN_clustered = X_numerics.copy()
            5  DBSCAN_clustered.loc[:,'Cluster'] = DBS_clustering.labels_
            6  # append labels to points
```

```
In [189]:   1  # With the following code, you will see that 5
            2  # clusters are created AND one outlier cluster (-1) was created.
            3  # You can also see that the size of clusters varies greatly,
            4  # and some only have a few observations.
            5  # Lastly, you see there are 18 outliers (-1).
            6  DBSCAN_clust_sizes = DBSCAN_clustered.groupby('Cluster').size().to_frame()
            7  DBSCAN_clust_sizes.columns = ["DBSCAN_size"]
            8  DBSCAN_clust_sizes
```
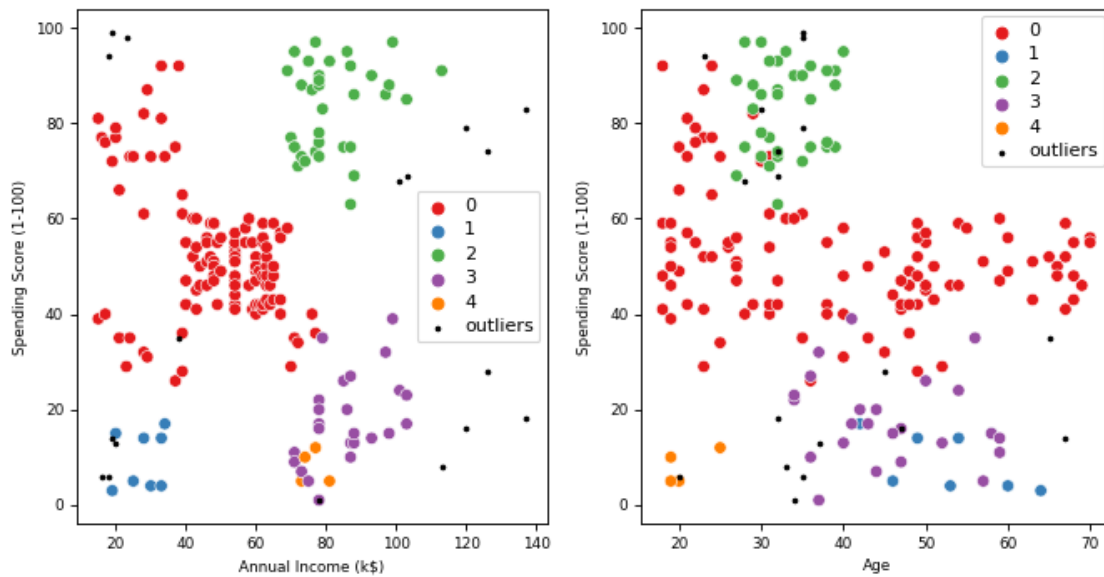
Out[189]:

|         | DBSCAN_size |
|---------|-------------|
| **Cluster** |         |
| **-1**  | 18          |
| **0**   | 112         |
| **1**   | 8           |
| **2**   | 34          |
| **3**   | 24          |
| **4**   | 4           |

```
1  # With this code, you will see that there are
2  # some outliers (small black dots).
3  #These points do not meet distance and minimum
4  # samples requirements to be recognized as a cluster.
5  outliers = DBSCAN_clustered[DBSCAN_clustered['Cluster']==-1]
6  fig2, (axes) = plt.subplots(1,2,figsize=(10,5))
7  sns.scatterplot(x='Annual Income (k$)', y='Spending Score (1-100)',
8  data=DBSCAN_clustered[DBSCAN_clustered['Cluster']!=-1],
9  hue='Cluster', ax=axes[0], palette='Set1', legend='full', s=45)
10 sns.scatterplot(x='Age', y='Spending Score (1-100)',
11 data=DBSCAN_clustered[DBSCAN_clustered['Cluster']!=-1],
12 hue='Cluster', palette='Set1', ax=axes[1], legend='full', s=45)
13 axes[0].scatter(outliers['Annual Income (k$)'], outliers['Spending Score (1-100)'],
14 s=5, label='outliers', c="k")
15 axes[1].scatter(outliers['Age'], outliers['Spending Score (1-100)'], s=5,
16 label='outliers', c="k")
17 axes[0].legend()
18 axes[1].legend()
19 plt.setp(axes[0].get_legend().get_texts(), fontsize='10')
20 plt.setp(axes[1].get_legend().get_texts(), fontsize='10')
21 plt.show()
```
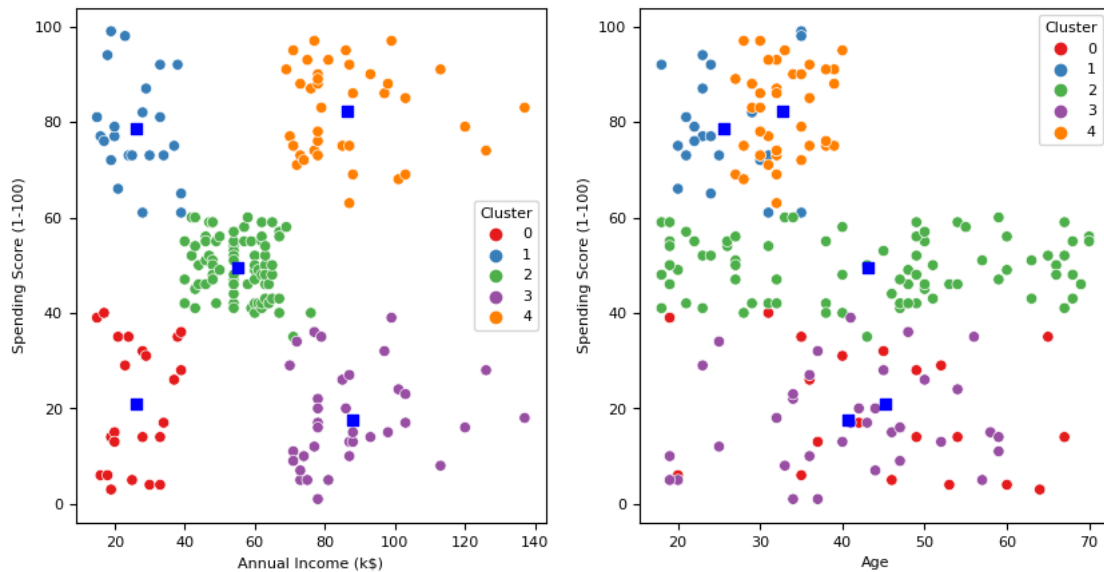
Figure 127



## STEP 7: Evaluate

Compare Algorithms

```
1  # K-Means with K=5
2  fig1.suptitle('K-Means', fontsize=16)
3  fig1
4
```

```
1  # K-Means with K=6
2  fig11.suptitle('K-Means', fontsize=16)
3  fig11
```
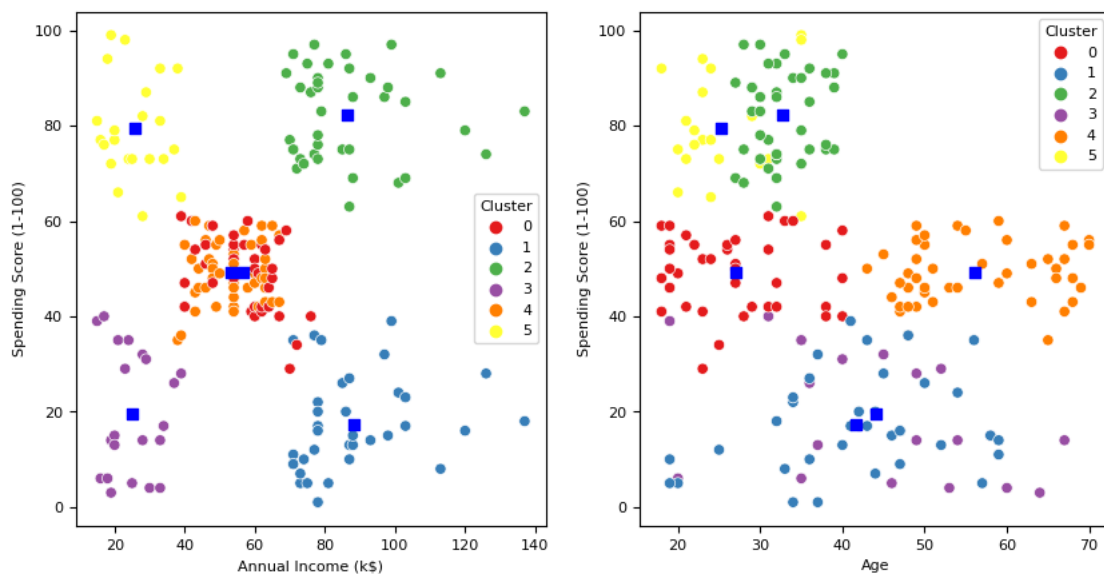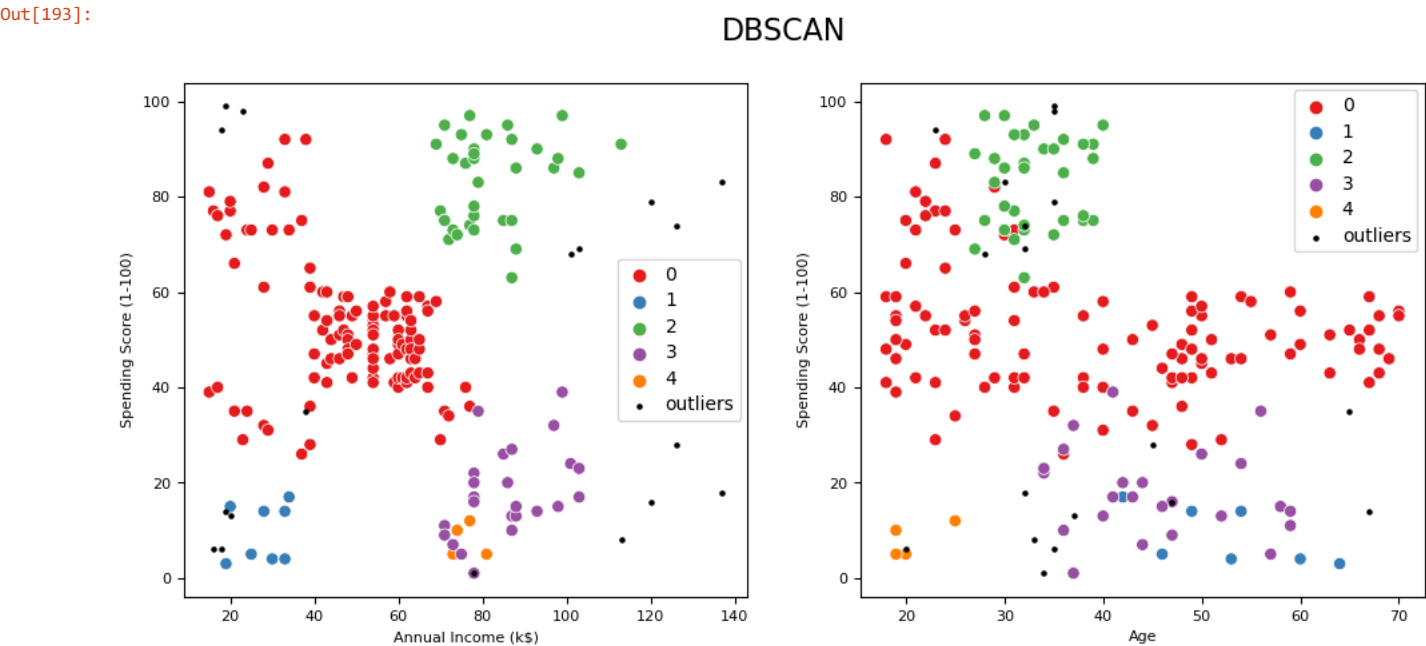
```
In [193]:  1  # DBSCAN with 5 clusters AND outliers
           2  fig2.suptitle('DBSCAN', fontsize=16)
           3  fig2
```

Out[193]:



```
In [194]:  1  clusters = pd.concat([KM6_clust_sizes, DBSCAN_clust_sizes],axis=1,sort=False)
           2  clusters
```

Out[194]:

| Cluster | KM_size | DBSCAN_size |
|---|---|---|
| 0 | 38.0 | 112.0 |
| 1 | 35.0 | 8.0 |
| 2 | 39.0 | 34.0 |
| 3 | 21.0 | 24.0 |
| 4 | 45.0 | 4.0 |
| 5 | 22.0 | NaN |
| -1 | NaN | 18.0 |

```
In [195]:  1  # Print the outliers
           2  print(DBSCAN_clustered[DBS_clustering.labels_==-1])
```

```
     Age  Annual Income (k$)  Spending Score (1-100)  Cluster
2    20                  16                       6       -1
6    35                  18                       6       -1
7    23                  18                      94       -1
10   67                  19                      14       -1
11   35                  19                      99       -1
14   37                  20                      13       -1
19   35                  23                      98       -1
40   65                  38                      35       -1
158  34                  78                       1       -1
187  28                 101                      68       -1
191  32                 103                      69       -1
192  33                 113                       8       -1
194  47                 120                      16       -1
195  35                 120                      79       -1
196  45                 126                      28       -1
197  32                 126                      74       -1
198  32                 137                      18       -1
199  30                 137                      83       -1
```
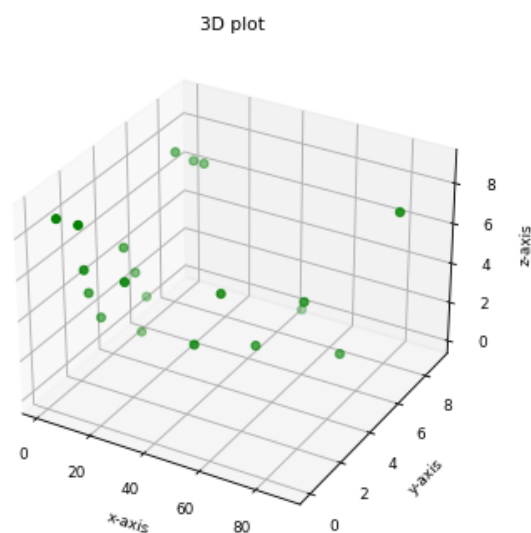
```
In [196]:   1  # creating 3d plot using matplotlib
            2  # in python
            3
            4  # for creating a responsive plot
            5  %matplotlib widget
            6
            7  # importing required libraries
            8  from mpl_toolkits.mplot3d import Axes3D
            9  import matplotlib.pyplot as plt
           10
           11  # creating random dataset
           12  xs = [14, 24, 43, 47, 54, 66, 74, 89, 12, 44, 1, 2, 3, 4, 5, 9, 8, 7, 6, 5]
           13
           14  ys = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 6, 3, 5, 2, 4, 1, 8, 7, 0, 5]
           15
           16  zs = [9, 6, 3, 5, 2, 4, 1, 8, 7, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 0]
           17
           18  # creating figure
           19  fig=plt.figure()
           20  ax=fig.add_subplot(111,projection='3d')
           21
           22  # creating the plot
           23  plot_geeks = ax.scatter(xs, ys, zs, color='green')
           24
           25  # setting title and labels
           26  ax.set_title("3D plot")
           27  ax.set_xlabel('x-axis')
           28  ax.set_ylabel('y-axis')
           29  ax.set_zlabel('z-axis')
           30
           31  # displaying the plot
           32  plt.show()
           33
```

Figure 119



3D plot

## Great Job! You are done.

**Real-world application questions to consider.**

Now that we have identified the outliers, now what?

- Why are they occurring?
- Where—and what—might the meaning be?

The answer could differ from organization to organization, but you can't simply ignore the data, regardless of the significance. In addition to the above questions, it is more important to know

- How outliers affect your testing efforts?
- How do I deal with the identified outliers?
- Should outliers be removed from the analysis?