

# Data Mining Techniques

## Spring Semester 2022-2023

### 1st Assignment – Work individually or in groups of 2

#### Assignment purpose

The purpose of this assignment is to learn the basic stages of the procedure followed to apply data mining techniques, namely: data preprocessing / cleaning, conversion, application of data mining techniques and evaluation. The implementation will be done in Python using tools / libraries: jupyter notebook, pandas and SciKit Learn.

#### Description

**Customer Profile Analysis** is a detailed analysis of a company's ideal customers. It helps companies understand their customers better and make the modification of its products according to the specific needs, behaviors and concerns of different types of customers easier. For example, instead of spending money for the promotion of a new product to each customer, a company can analyze which customer category is most likely to buy that product and then promote it only to that particular category.

## About data

### Information about Customers

- **ID:** Customer's unique identifier
- **Year\_Birth:** Customer's birth year
- **Education:** Customer's education level
- **Marital\_Status:** Customer's marital status
- **Income:** Customer's yearly household income
- **Kidhome:** Number of children in customer's household
- **Teenhome:** Number of teenagers in customer's household
- **Dt\_Customer:** Date of customer's enrollment with the company
- **Recency:** Number of days since customer's last purchase
- **Complain:** 1 if the customer complained in the last 2 years, 0 otherwise.

### Products (money spent in two years)

- **MntWines:** Amount spent on wine in last 2 years
- **MntFruits:** Amount spent on fruits in last 2 years
- **MntMeatProducts:** Amount spent on meat in last 2 years
- **MntFishProducts:** Amount spent on fish in last 2 years
- **MntSweetProducts:** Amount spent on sweets in last 2 years
- **MntGoldProds:** Amount spent on gold in last 2 years.

### Promotion

- **NumDealsPurchases:** Number of purchases made with a discount
- **AcceptedCmp1:** 1 if customer accepted the offer in the 1st campaign, 0 otherwise
- **AcceptedCmp2:** 1 if customer accepted the offer in the 2nd campaign, 0 otherwise
- **AcceptedCmp3:** 1 if customer accepted the offer in the 3rd campaign, 0 otherwise
- **AcceptedCmp4:** 1 if customer accepted the offer in the 4th campaign, 0 otherwise
- **AcceptedCmp5:** 1 if customer accepted the offer in the 5th campaign, 0 otherwise
- **Response:** 1 if customer accepted the offer in the last campaign, 0 otherwise

### Source

- **NumWebPurchases:** Number of purchases made through the company's website
- **NumCatalogPurchases:** Number of purchases made using a catalogue
- **NumStorePurchases:** Number of purchases made directly in stores
- **NumWebVisitsMonth:** Number of visits to company's website in the last month

## Questions

1. **Preprocessing / Cleaning:** Check if there are missing values in the data and handle them accordingly, convert columns about dates to DateTime objects and check for any dtype: object attributes that can be encoded / converted to numeric values **(5%)**.
2. Print the **unique values** of categorical attributes **Marital\_Status** and **Education** to get a more clear picture of the data. Change values [Alone, Absurd, YOLO] of Marital\_Status to 'Single'. Use any type of graph to show the number of values in each category **(5%)**.
3. **Create new attributes (10%):**
  - A. Create an attribute "**Customer\_For**" that represents the number of days passed from the day the customers making purchases from the store in relation to the last recorded date (Recency).
  - B. Extract the age "**Age**" of a customer based on "**Year\_Birth**" which indicates the respective person's year of birth.
  - C. Create another attribute, "**Spent**", which indicates the customer's total money spent in all categories in two years.
  - D. Create an attribute "**Children**" that shows the number of children in a household, i.e. kids and teenagers.
  - E. To get a better understanding of a household, create an attribute named "**Family\_Size**" that shows the total number of people in a household.
  - F. Create an attribute "**Is\_Parent**" that indicates whether or not a customer is a parent.
  - G. Create another attribute, "**Living\_With**", using "**Marital\_Status**" to extract the living situation of couples. Specifically, this attribute must have two values. "Partner" and "Alone".
  - H. Create column "**Age\_Group**" using column "**Age**", to group ages according to the following values: "21-30", "31-40", "41-50", "51-60", "61-70", "71-80", ">80".
4. Check for any **outliers** in the attributes and delete them from the data **(5%)**.
5. Then, examine the **correlation of the attributes** using a heatmap (excluding categorical attributes at this point) **(5%)**.

6. Answer the following questions with **graphs** – choose 10 of the following **(20%)**
1. In which **Marital\_Status** category does the largest percentage of the customers belong?
  2. How many customers have filled a **Complain**?
  3. What is the relationship between the number of purchases **Spent** and **Marital\_Status**?
  4. What is the relationship between the number of purchases **Spent** and **KidHome** and **Family\_Size** attributes?
  5. What is the relationship between **Age\_Group** and **Spent** attributes?
  6. What is the relationship between **Income** and **Spent** attributes?
  7. What is the relationship between **Education** and **Income** attributes?
  8. What is the relationship between **Income** and **Family\_Size** attributes?
  9. What is the relationship between **Income** and **KidHome** attributes?
  10. What is the relationship between **Income** and **Living\_With** attributes?
  11. What is the relationship between **Income** and **Spent** attributes?
  12. What is the relationship between **NumWebPurchases** and **NumWebVisitsMonth** attributes?
  13. What is the percentage of customers that has accepted all deals from the store?
  14. Draw a histogram for **NumDealsPurchases** column.
  15. Draw a histogram for **Income** column.
  16. Draw a histogram for **KidHome** column.
  17. Draw a histogram for **Family\_Size** column.
  18. Do customers with a masters degree spent more money on wine?

7. **Principal Component Analysis (PCA) (25%):**

In this problem, there are lots of factors on which a classification is made. These factors are basic attributes of features. The greater the number of attributes, the harder the problem is. Many of these attributes are correlated and thus, redundant. This is the reason why you are going to preform **dimension reduction** to specific attributes. Dimension reduction is the process of reducing the number of random variables that are under examination and leads to obtaining of a set of major variables.

After the addition of attributes done in the previous questions, the dataset variables that are categorical and not numeric, are the following ['**Education**', '**Marital\_Status**', '**Living\_With**']. Use **LabelEncoder()** for these variables to convert them to numeric data (this process is called **one hot encoding**).

Then, create a copy of the dataframe, which contains all numerical columns and delete columns related to offers and marketing campaigns, i.e. ['**AcceptedCmp1**', '**AcceptedCmp2**', '**AcceptedCmp3**', '**AcceptedCmp4**', '**AcceptedCmp5**', '**Complain**', '**Response**'].

In this way, the data obtained contains attributes of various dimensions and variances. Different variances in attributes negatively affect the modeling of the dataset. The solution is to perform *Standardization*, so that each column / attribute / variable has  $\mu = 0$  and  $\sigma = 1$ .

Finally, use Principal Component Analysis (PCA) to reduce the number of dimensions to  $n\_components = 3$ . Draw the (three dimensional) graph of the result.

## 8. Clustering (25%)

### Steps

- Elbow method to determine the number of clusters needed.
- Agglomerative and K-Means clustering
- Show the formed clusters with a graph (e.g. scatter plot).

## 9. Customer Profile Analysis (Bonus)



Try to sketch the profile of the formed clusters through graphs to arrive to a conclusion regarding who is the most "important" customer and who needs the most attention from the store's marketing team.

To achieve this, provide some of the attributes that are indicative of the personal

characteristics of a customer depending on the cluster they belong to (e.g. **Age**, **Is\_Parent**, **Family\_Size** etc.). Finally, collect the main characteristics of each cluster.

e.g.

<b>Cluster 0:</b> They spend the least They have the lowest income They have teenagers at home They are older	<b>Cluster 1:</b> They spend the most They have a large income Most of them aren't parents They have taken part in all 6 marketing campaigns
---	--