# ESILV- Python for Data Analysis - Project

**Written by Jihane Oul Ali and Nada Matrouf**

# Content

- Data Introduction
- Data Visualization
- Data Preparation
- Modelisation and Validation
- Conclusion
- References

# 1. Data Introduction

**Abstract**

The data represents the estimation of obesity levels in individuals from Mexico, Peru and Columbia. The estimation is based on their eating habits and physical condition. The data contains 17 attributed with 2111 instances.

The target is labeled with the following values :

- *Insufficient Weight*

- *Normal Weight*

- *Overweight Level I*

- *Overweight Level II*

- *Obesity Type I*

- *Obesity Type II*

- *Obesity Type III*

# 1. Data Introduction

**Link:**

*https://archive.ics.uci.edu/ml/datasets/Estimation+of+obesity+levels+based+on+eating+habits+and
+physical+condition+*

| Data Set Characteristics: | Multivariate | Number of Instances: | 2111 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Integer | Number of Attributes: | 17 | Date Donated | 2019-08-27 |
| Associated Tasks: | Classification, Regression, Clustering | Missing Values? | N/A | Number of Web Hits: | 20563 |

# 1. Data Introduction

The attributes of the dataset and their input are based on the following:

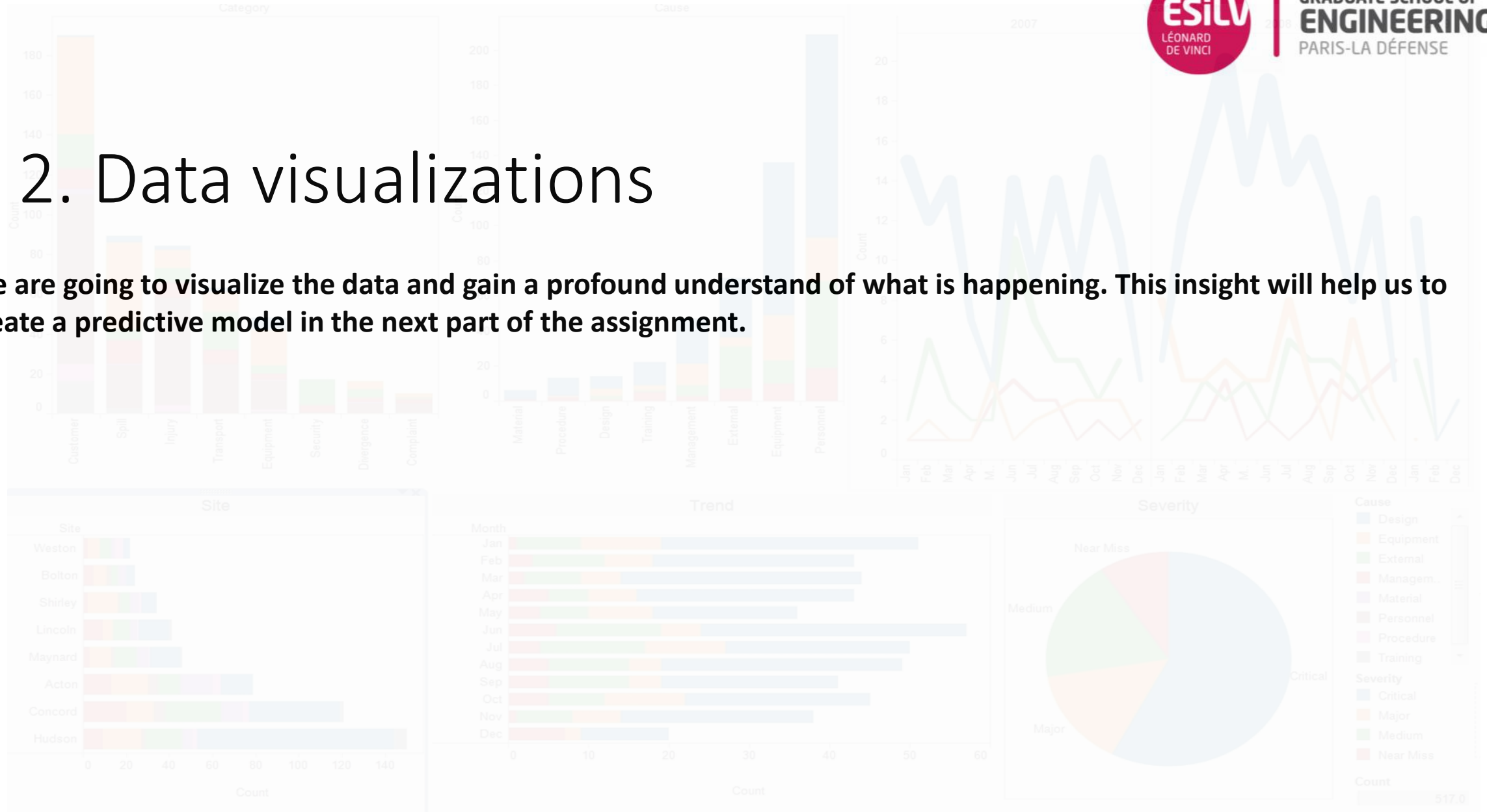| | | | |
|---|---|---|---|
| **Gender** | *Male \|Female* | **Food consumption between meals** | *No \| Sometimes \| Frequently \| Always* |
| **Age** | *Numeric Value* | **Smoking** | *Boolean* |
| **Height** | *Num. Val. In meters* | **Water consumption** | *Less than 1L\| Between 1 and 2L\| More than 2L* |
| **Weight** | *Num. Val. In Kg* | **Are calories monitored daily** | *Yes \| No* |
| **Family member was or is overweight** | *Boolean* | **Frequency of physical activity** | *I do not have \| 1 or 2 days\| 2 or 4 days \| 4 or 5 days* |
| **Frequent consumption of high caloric food** | *Boolean* | **Time spent on technological devices** | *0-2 hours \| 3-5 hours\| more than 5 hours* |
| **Frequency of vegetable consumption** | *Never \| Sometimes \| Always* | **Frequency of alcohol consumption** | *I do not drink\| sometimes\| frequently \|always* |
| **Number of daily meals** | *Between 1 and 2 \| Three\| More than three* | **Type of transportation used** | *Auto\| motorbike \|bike \|public transport \| walking* |

# 1. Data Introduction

**The goal :**

*Create a prediction model able to predict in which of the obesity classification labels an individual fits.*
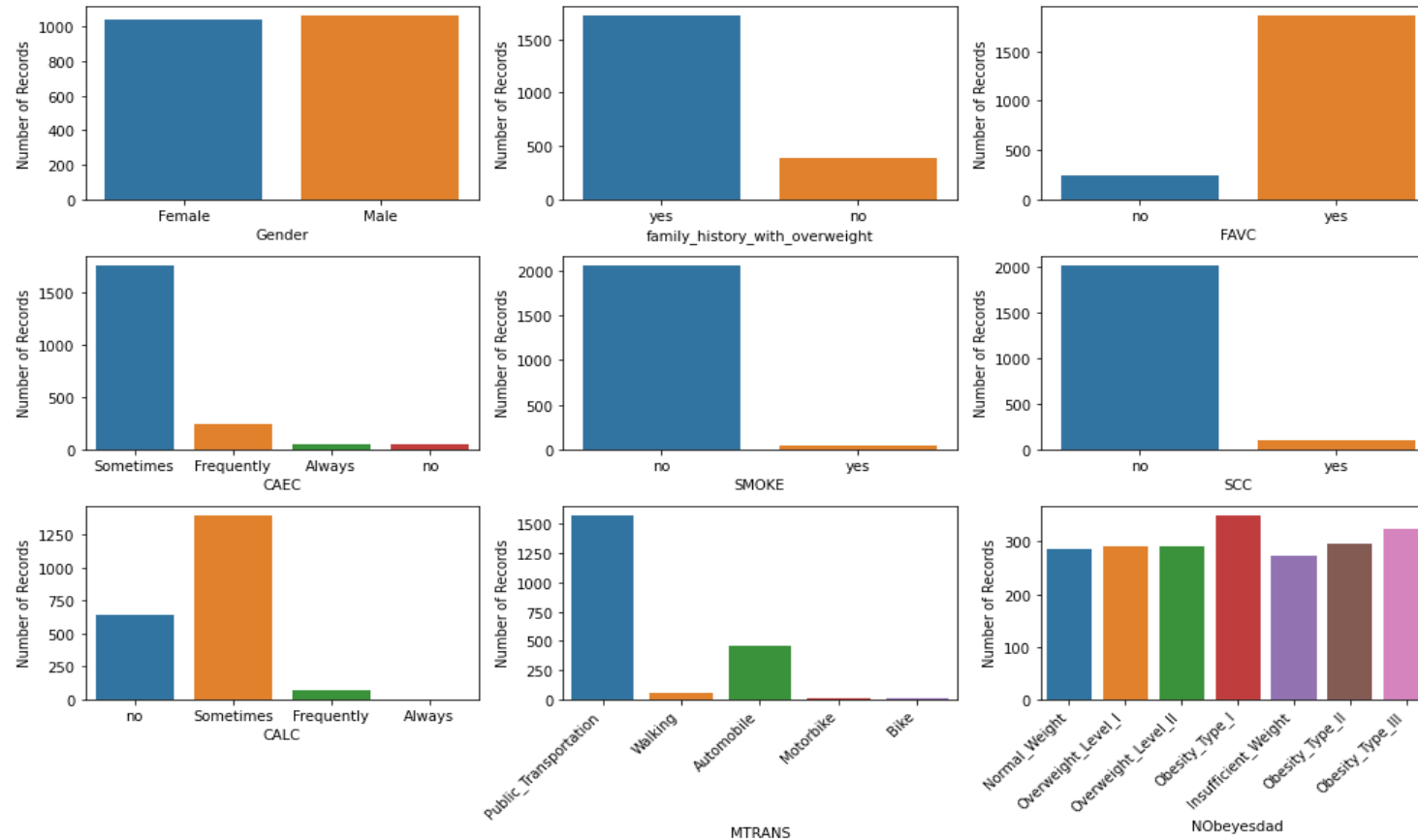
However, since the height and weight of an individual are also attribute, the classification can easily be done using the BMI index of a person. In order to analyze the data accordingly a model will also be created while removing the two attributes able to determine the BMI.

# 2. Data visualizations

**We are going to visualize the data and gain a profound understand of what is happening. This insight will help us to create a predictive model in the next part of the assignment.**
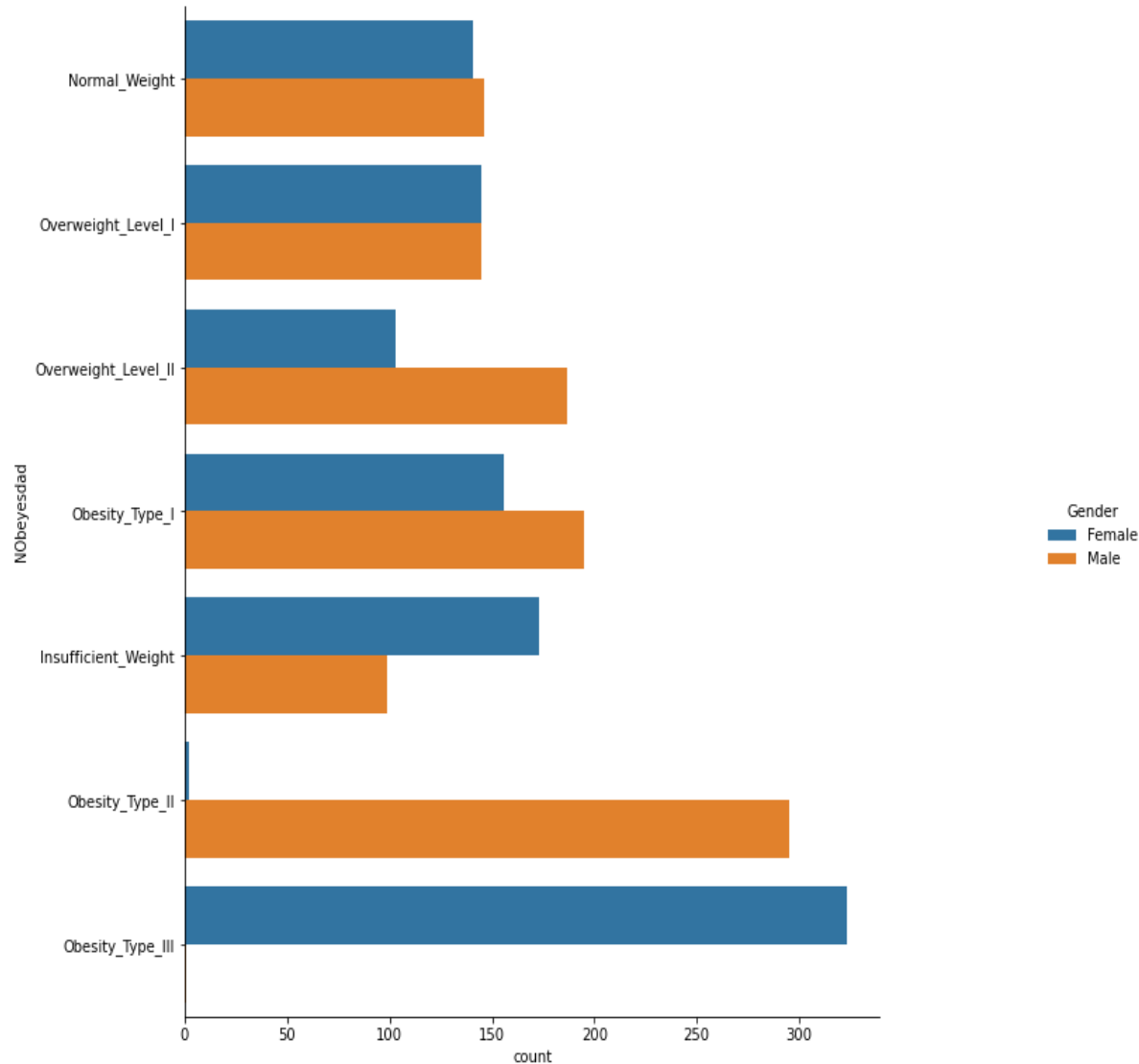
Attribute inputs

First, we will visualize the data that have various categories. This way all the different Input variables can be studied. From the histograms we can see that:
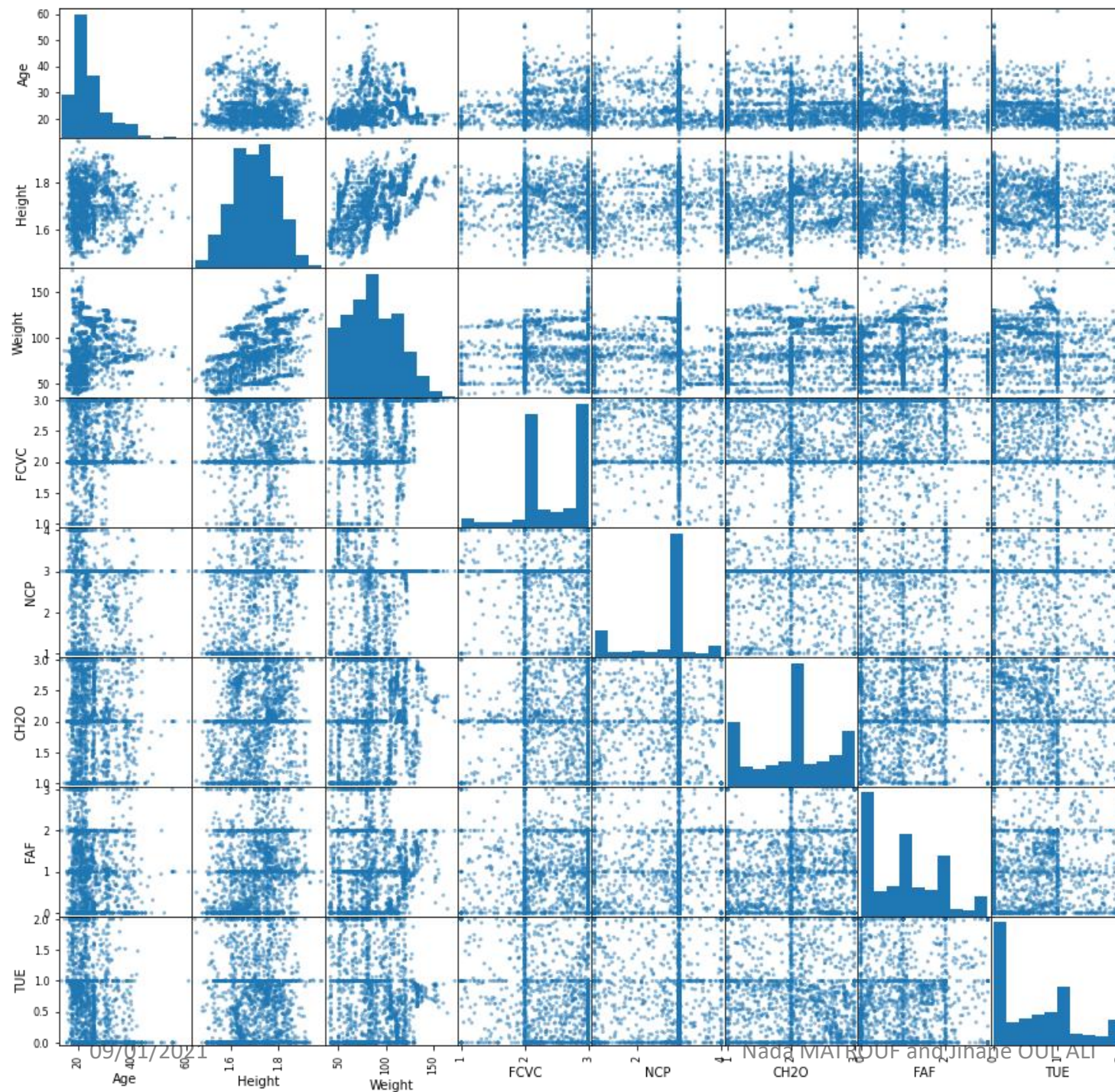
- There is an even number of female and male participant.

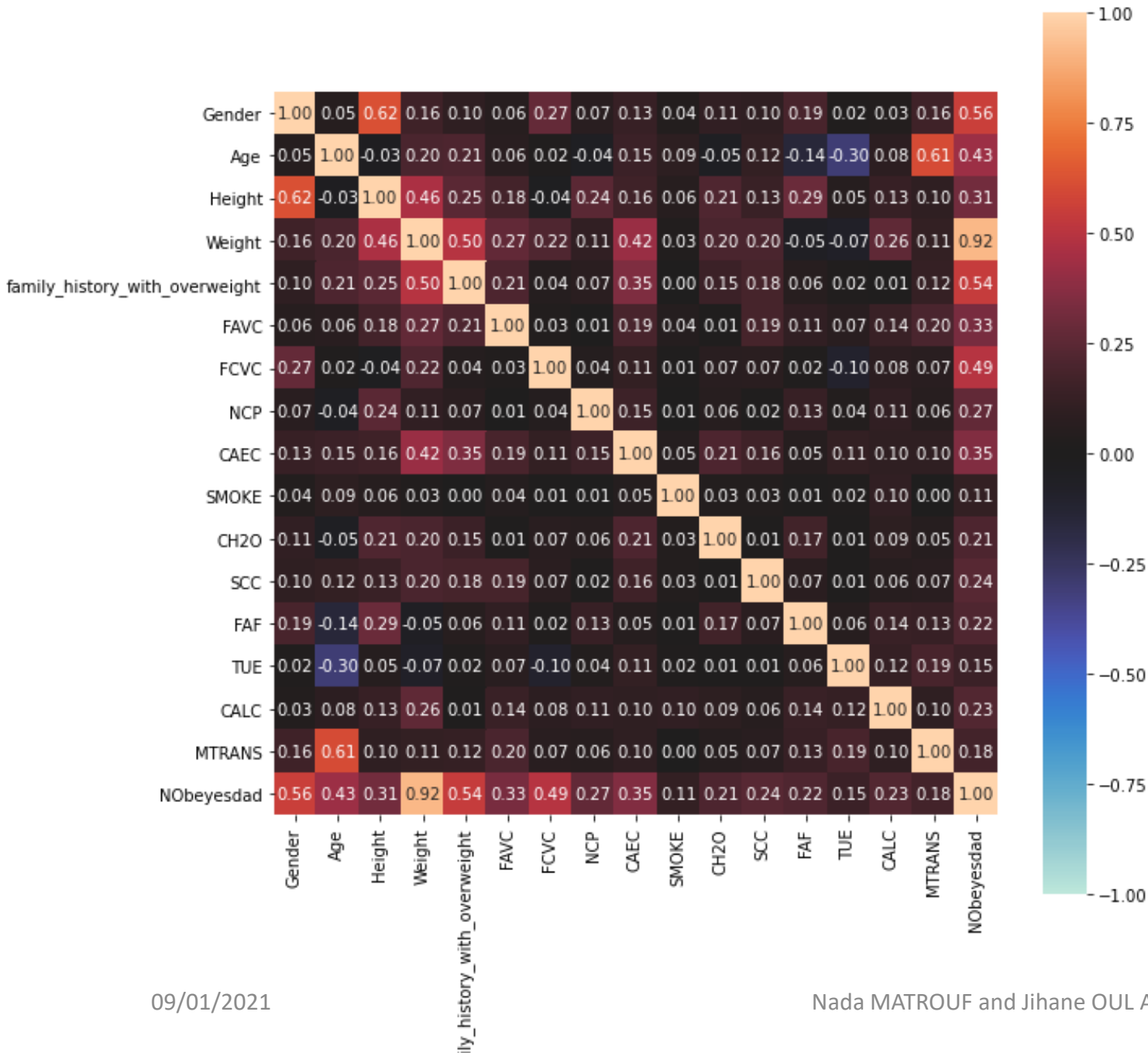The rest of the attributes show a clear distinction between each category.

Even though the number female and male participant is the same, a clear distinction can be made in the target values for almost all of the classes except the normal weight class and the first type of overweight class. This indicates that the gender plays a very important role in assigning the labels to each individual.

Nada MATROUF and Jihane OUL ALI

The following plot shows a scatter plot of each of the 15 attributes.

A heat map is also created. The heatmap shows us the correlation between each attribute and therefore its importance in the prediction model.

As previously explained, the correlation between the target attribute and weight seems very high. Which agrees with our previously made hypothesis that the results heavenly rely on the BMI.

However, the heat map also shows a high correlation between the Gender and the target attribute NObeyesdad. The next important attributes are :
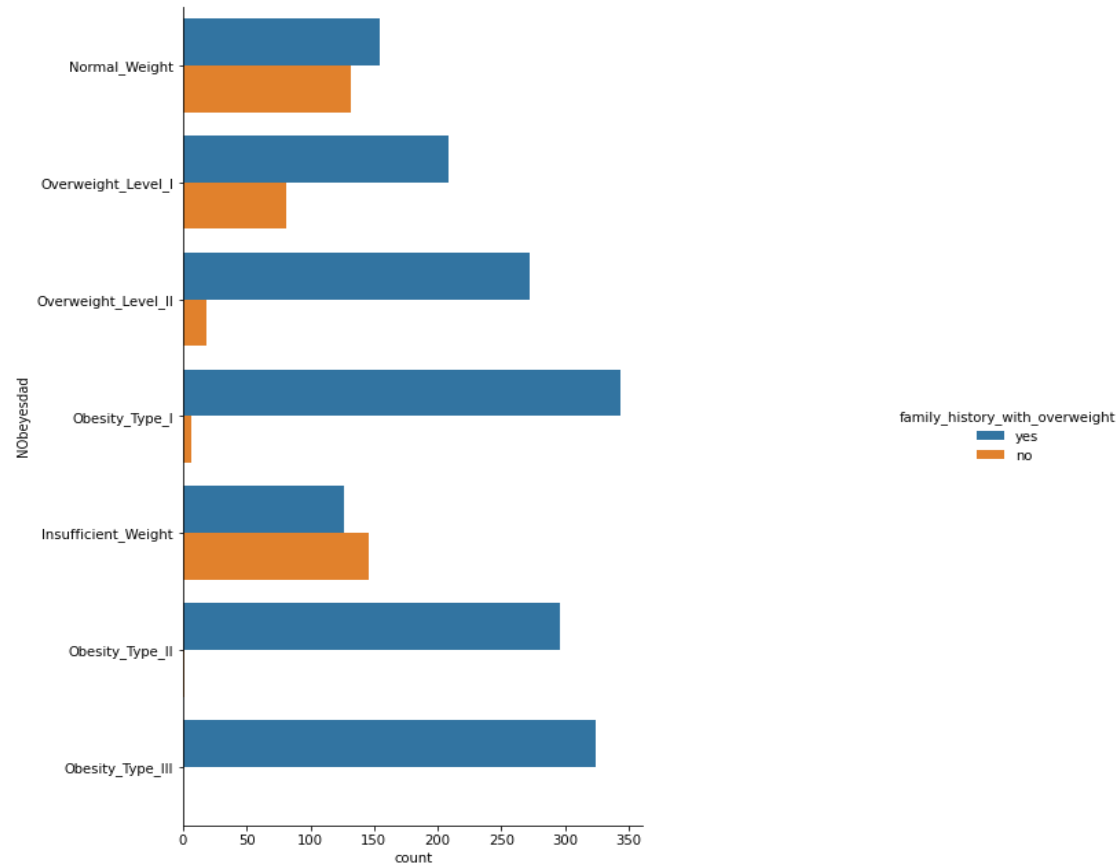
*Family history with, overweight,FCVC,AGE,CAEC,FAVC* in that order.

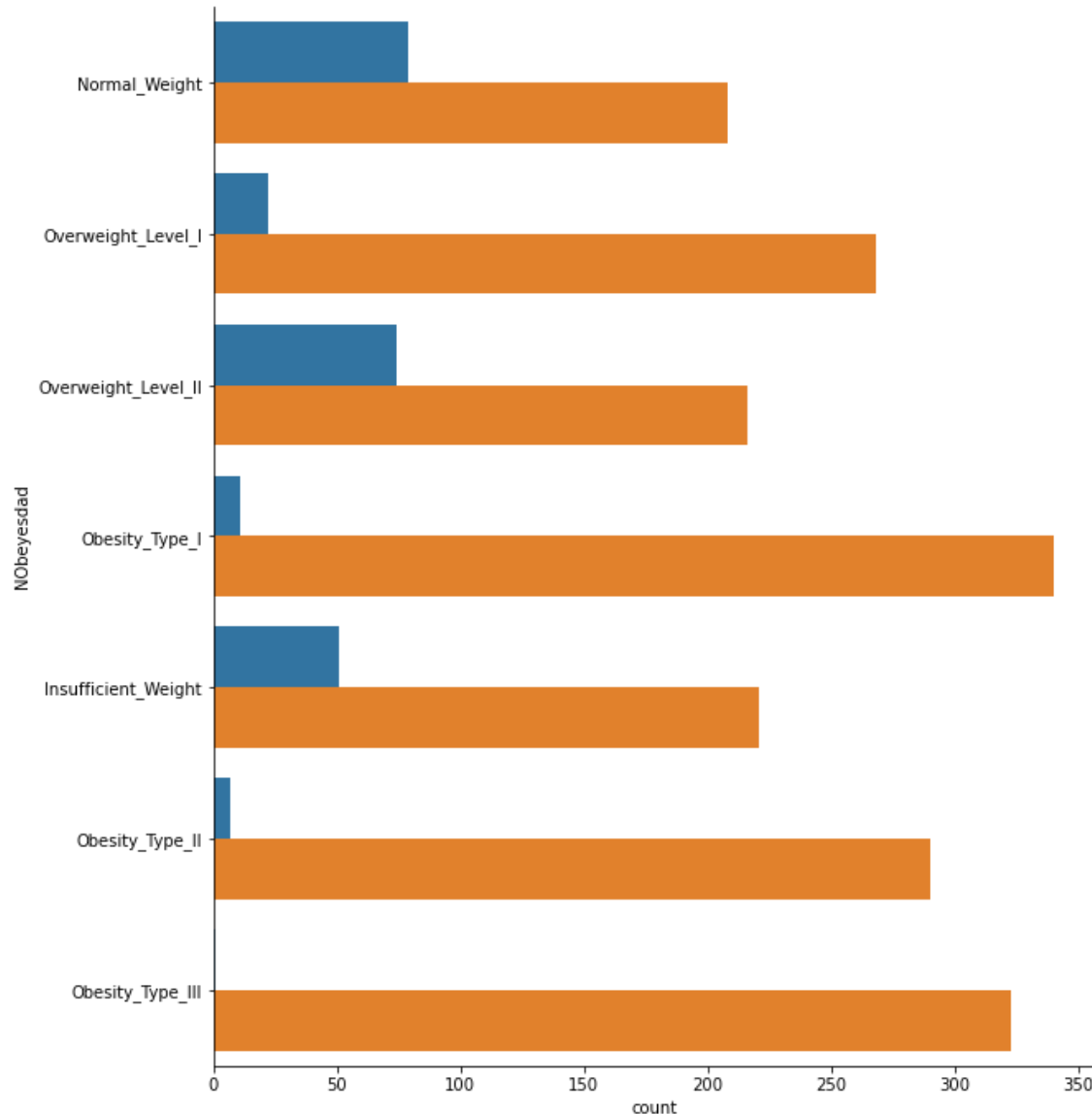Some of these attributes will be looked at closer during the visualization phase to gain more insight.

# Family history with overweight

First,we will start with the attribute showcasing the family history with overweight. The input values for this attribute is Boolean.

We can clearly see that although for all the target except obesity type 2 and obesity type 3, there are some individuals suffering even though no history of overweight family members is detected. However, it can't be denied that the influence of family history is huge on all types of obesity and overweight type 1 and 2. This agrees with the many studies performed previously on the link of overweight family history and severe onset of obeisant. Below links of some of the studies performed:

# Frequency consumption of high caloric food



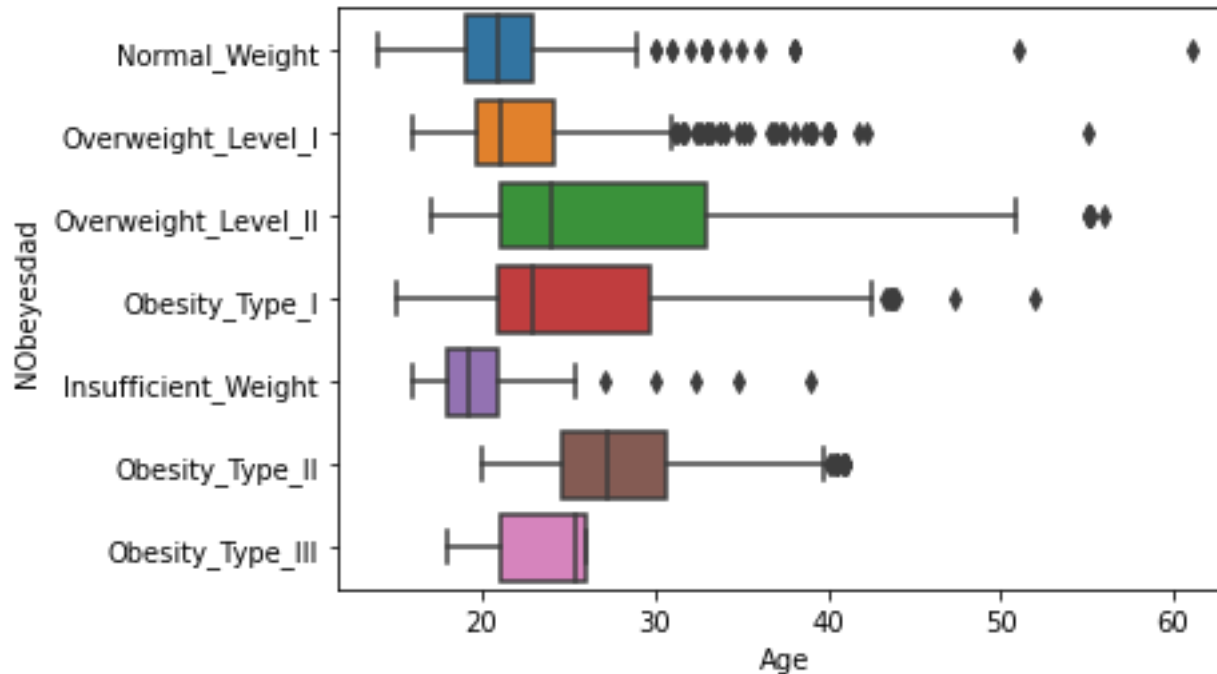The next step is to determine how the frequent consumption of food high in calories influences the various types of target labels:

Here we can see that most of the people answered with yes which mean the high amount of calories does not necessarily influence the output. However, it can be noticed that the people with a normal weight do have a tendency to eat less high caloric meals on a daily basis.
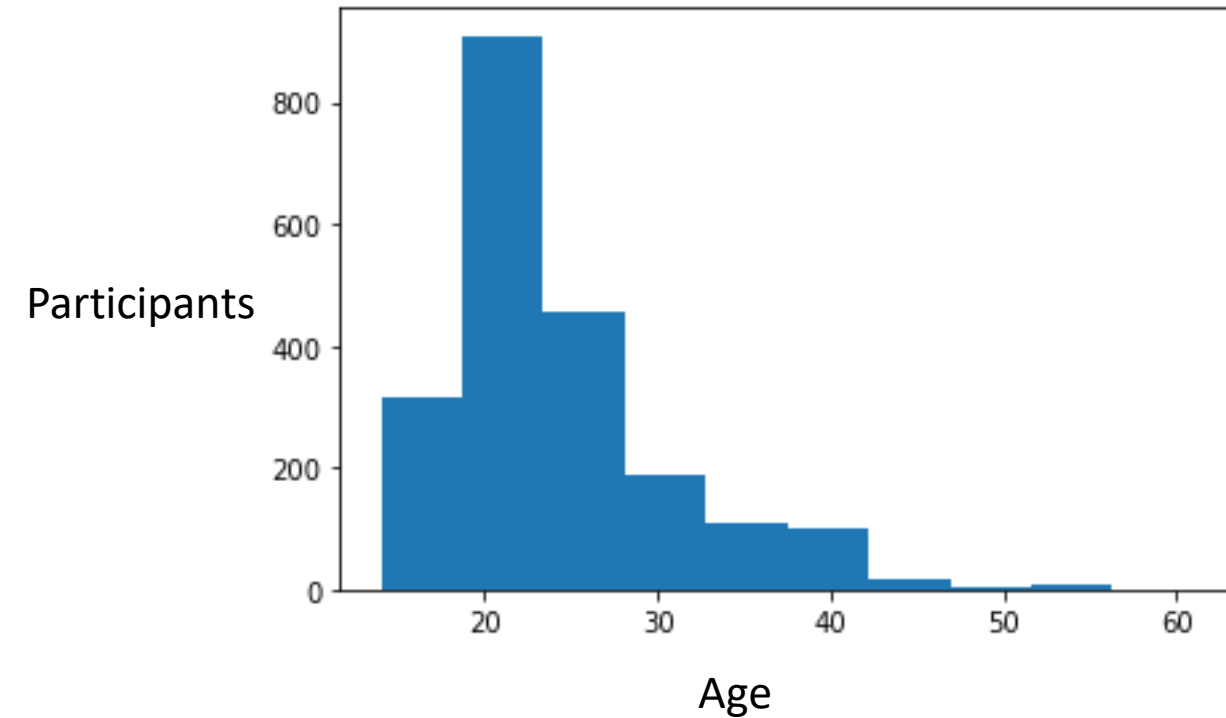
| family_history_with_overweight | no | | yes | |
|---|---|---|---|---|
| FAVC | no | yes | no | yes |
| NObeyesdad | | | | |
| Insufficient_Weight | 46.0 | 100.0 | 5.0 | 121.0 |
| Normal_Weight | 36.0 | 96.0 | 43.0 | 112.0 |
| Obesity_Type_I | 1.0 | 6.0 | 10.0 | 334.0 |
| Obesity_Type_II | 1.0 | NaN | 6.0 | 290.0 |
| Obesity_Type_III | NaN | NaN | 1.0 | 323.0 |
| Overweight_Level_I | 7.0 | 74.0 | 15.0 | 194.0 |
| Overweight_Level_II | 8.0 | 10.0 | 66.0 | 206.0 |

Here we can see that most of the people that have a family members with history of obesity also tend to consume high calorie meals on a daily bases. The majority of these people have type 1 obesity closely followed by type 3 obesity. This is a very important remark since it tends to also show the link between having a family member that suffered or is suffering from obesity and the tendency to eat more.

# Age

Next we will look at the age. The following boxplots will show the average age for each target class.



We can see that most of the targets have an average age between 20 and 30. This sounds counter intuitive since most of the people at that age are at their fittest. Therefore, we will plot a graph showcasing the age of the individuals.

Participants

Age

The plot shows that the majority of the participants is between 20 and 30. This explains our previously found conclusions. However, the boxplots are not completely useless. The boxplots showcasing an average that leans toward the 30 indicates a very strong influence for the people that are above 30. Since this group of participants is small the numbers must be high. These include obesity type 2 and type 3.

# 3. Data Preparation

As explained previously, we will create two models one with the height and weight and one without. In addition, the different strings in the attributes will be transformed to categories. (Other encoding processes are possible with different libraries however this seemed to be the most straightforward)

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2111 entries, 0 to 2110
Data columns (total 17 columns):
 #   Column                          Non-Null Count  Dtype
---  ------                          --------------  -----
 0   Gender                          2111 non-null   category
 1   Age                             2111 non-null   float64
 2   Height                          2111 non-null   float64
 3   Weight                          2111 non-null   float64
 4   family_history_with_overweight  2111 non-null   category
 5   FAVC                            2111 non-null   category
 6   FCVC                            2111 non-null   category
 7   NCP                             2111 non-null   category
 8   CAEC                            2111 non-null   category
 9   SMOKE                           2111 non-null   category
 10  CH2O                            2111 non-null   float64
 11  SCC                             2111 non-null   category
 12  FAF                             2111 non-null   category
 13  TUE                             2111 non-null   float64
 14  CALC                            2111 non-null   category
 15  MTRANS                          2111 non-null   category
 16  NObeyesdad                      2111 non-null   category
dtypes: category(12), float64(5)
memory usage: 235.5 KB
```

We will also use encoding for the CAEC and CALC perimeters since the degree here is important.

# 4. Modelisation and Validation

In order to build a prediction model many machine learning algorithms can be applied. For this specific application the goal is to create a prediction model able to classify the individual based on the given attribute. This leaves us with the many options for classification algorithms. For this approach three types will be considered:

- **Decision tree model** : In this model the algorithm creates many subtrees depending on the attribute input.

- **Random Forest model** : This is an extension of the decision tree model. The Random Forest Algorithm combines the output of multiple (randomly created) Decision Trees to generate the final output.

- **KNN** : KNN works by finding the distance between data points by selecting a number K of groups and classifying them by distance.

# Dataset without the height and the weight

```
##########Decision tree Model##########
Precision of the DT Model: 0.8474 %
Recall of the DT Model: 0.8474 %
f1-score of the DT Model: 0.2615 %
##########Random Forest Model##########
Precision of the DT Model: 0.9726 %
Recall of the DT Model: 0.9726 %
f1-score of the DT Model: 0.4909 %
##########KNN Model##########
Precision of the DT Model: 0.6282 %
Recall of the DT Model: 0.6282 %
f1-score of the DT Model: 0.1811 %
```

From these results, we can conclude that the Random Forest Model proved to be the best model for our data set with an accuracy of 97%.

# Dataset with the height and the weight

Next we will look at the dataset while including the height and the weight.

```
##########Decision tree Model##########
Precision of the DT Model: 0.8121 %
Recall of the DT Model: 0.8121 %
f1-score of the DT Model: 0.2553 %
##########Random Forest Model##########
Precision of the RF Model: 0.4207 %
Recall of the RF Model: 0.4207 %
f1-score of the RF Model: 0.2031 %
##########KNN Model##########
Precision of the KNN Model: 0.0000 %
Recall of the KNN Model: 0.0000 %
f1-score of the KNN Model: 0.0000 %
```

# 5. Conclusion

Although for the last models the height and weight were added, the predictions were not more accurate. The Random Forest Model scored the highest with an accuracy of 93%. However, the Random Forest Model accuracy made without the two attributes : height and weight is still higher (97%). Therefore, the Random Forest Model will be chosen for the rest of this assignment.

# References

Special thanks to

Fabio Mendoza Palechor and Alexis de la Hoz Manotas

For making the dataset available