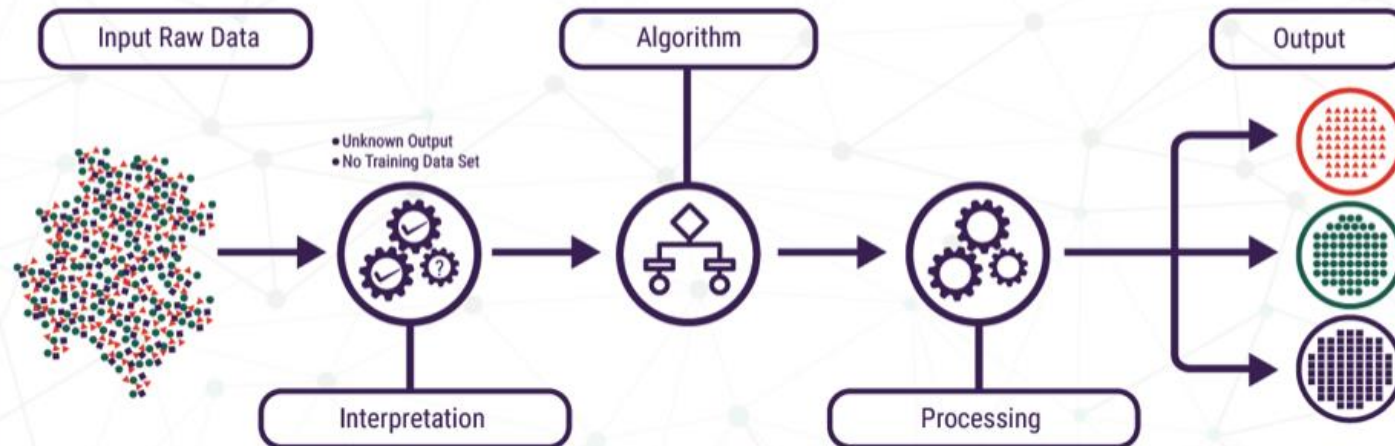# PYTHON FOR DATA ANALYSIS PROJECT

_

## Antoine PHISSAMAY Théo PHILIPPE



## 1. Learning from the dataset

First, in order to face a data analysis problem, we need to understand on what we work. Obviously, we need to open the dataset and see what we are facing. We observe 10 columns, and 5472 rows.

The dataset is about a project of making a decision tree more efficient. In fact, the dataset is composed by block charactristics, there are five different type of blocks :

- Vertical line
- Horrizontal line
- Text
- Pictures
- Graphics

Thus, the aim of the study is to categorize whether if rows in the dataset are data from one type of block or another. To deal with that we need to develop a machin learning model, here we don't have a training set with target values so we have to do an unsupervised learning.

*Dataset before preparation*

## 2. Data Cleaning and Preparation

The dataset is filled with columns aligned numbers, hence we can't use any function with python for data analysis. Before cleaning the data we need to make the file readable by cleaning and other functions.

We create a program which delete all the spaces and put each data in lists which separate each data by a coma. Then we also make a csv file with the data.
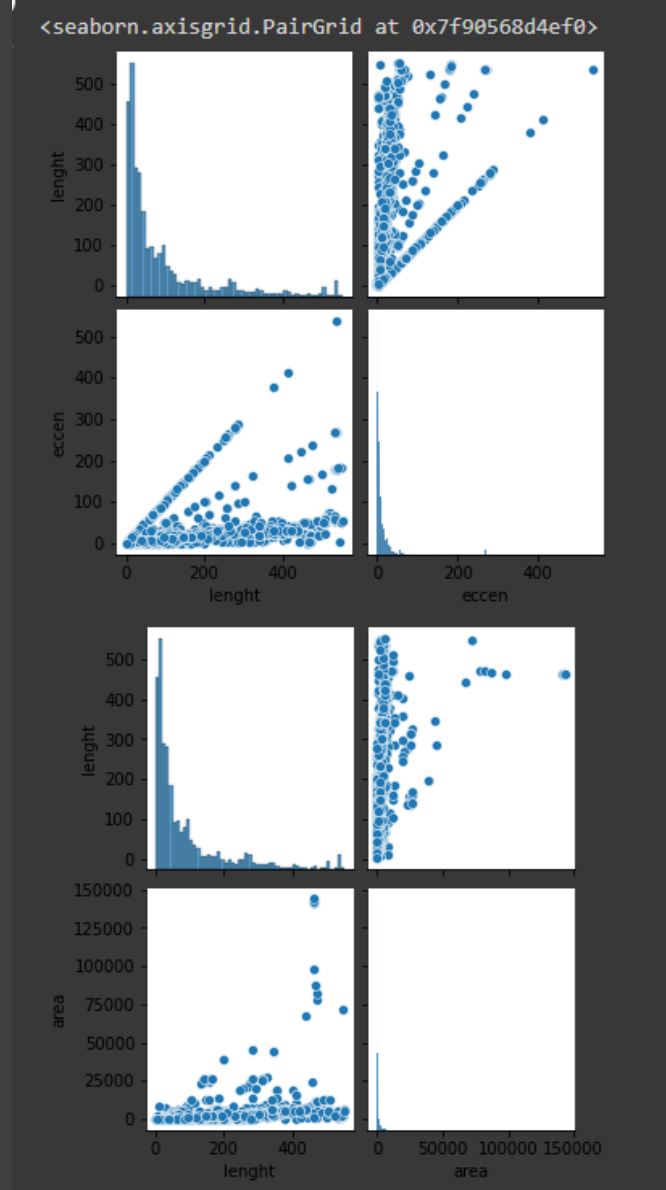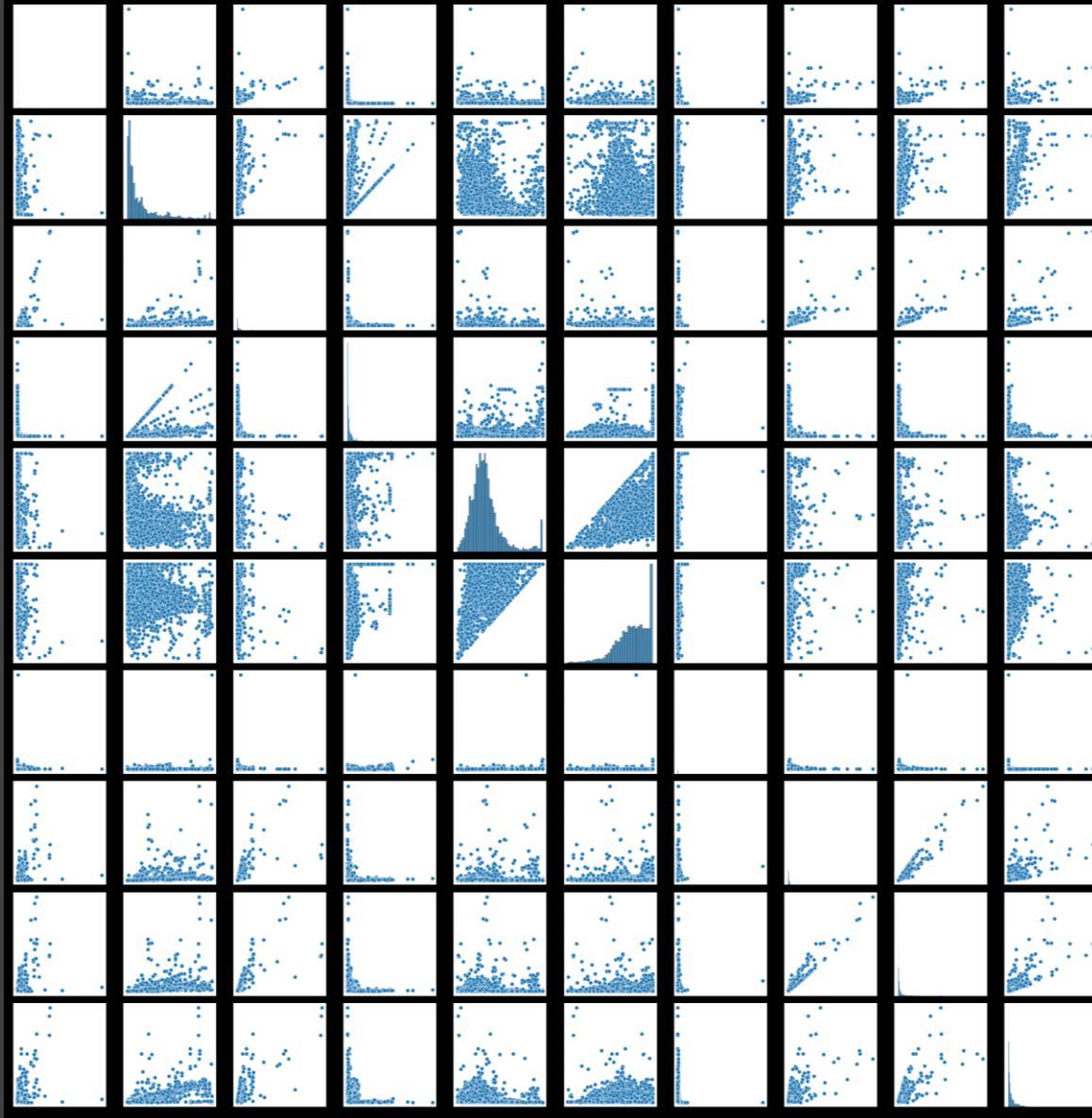
However, like it is specified in the information file attached with the dataset, there is no missing value, so we don't need to clean the dataset by deleting rows where values are missing. We can begin to work on the data.

*Dataset after preparation*

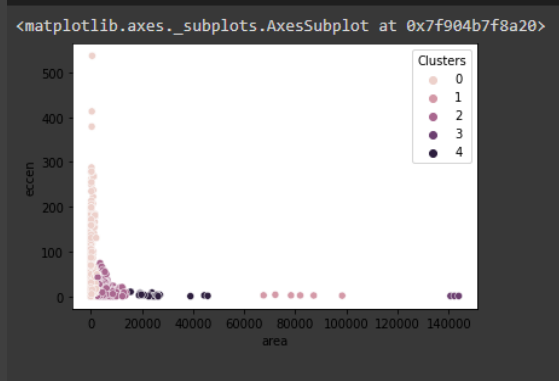| | height | lenght | area | eccen | p_black | p_and | mean_tr | blackpix | blackand | wb_trans |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 5 | 7 | 35 | 1.400 | 0.400 | 0.657 | 2.33 | 14 | 23 | 6 |
| 1 | 6 | 7 | 42 | 1.167 | 0.429 | 0.881 | 3.60 | 18 | 37 | 5 |
| 2 | 6 | 18 | 108 | 3.000 | 0.287 | 0.741 | 4.43 | 31 | 80 | 7 |
| 3 | 5 | 7 | 35 | 1.400 | 0.371 | 0.743 | 4.33 | 13 | 26 | 3 |
| 4 | 6 | 3 | 18 | 0.500 | 0.500 | 0.944 | 2.25 | 9 | 17 | 4 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 5467 | 1 | 16 | 16 | 16.000 | 1.000 | 1.000 | 16.00 | 16 | 16 | 1 |
| 5468 | 4 | 524 | 2096 | 131.000 | 0.542 | 0.603 | 40.57 | 1136 | 1264 | 28 |
| 5469 | 7 | 4 | 28 | 0.571 | 0.714 | 0.929 | 10.00 | 20 | 26 | 2 |
| 5470 | 6 | 95 | 570 | 15.833 | 0.300 | 0.911 | 1.64 | 171 | 519 | 104 |
| 5471 | 7 | 41 | 287 | 5.857 | 0.213 | 0.801 | 1.36 | 61 | 230 | 45 |

5472 rows × 10 columns

>> *'Panda.read_csv()'*

# 3. Data Visualization

Before we make a model, we need to see data repartition and correlations between parameters.
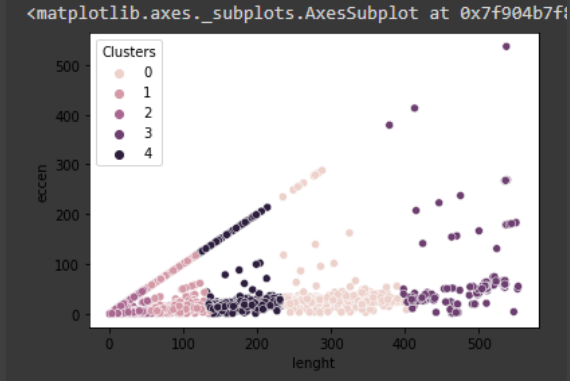
Displays of correlations shows that a few parameters are not much correlate.

It seems like some features are more correlate than others. We observe relations in order to build our model.

*Correlation between area and eccen*



*Correlation between lenght and eccen*



*Correlation between area and lenght*

# 4. Building of the model

As said before, this problem is about predicting a category.

Moreover, the data is considered as labeled data as fa as we need to categorize it. The Machin Learning algorithm to use here is a clustering algorithm.
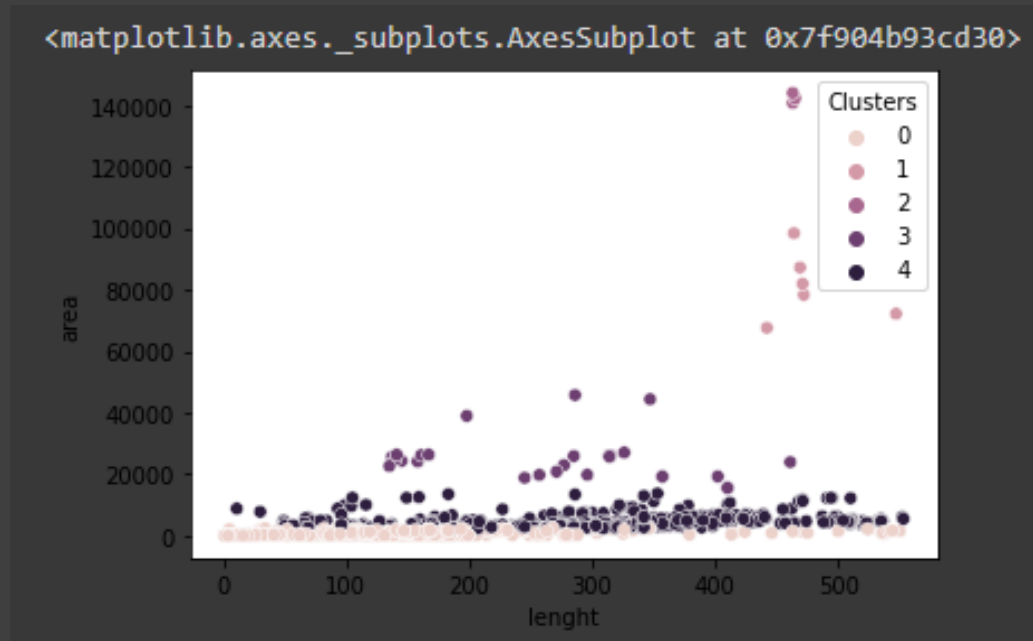
The number of categories is known so we can operate a K-Means algorithm.

Using SciKit Learn we build a model with the K-Means program.

Furthermore, correlation graphs before, expose which pair of features are the less correlate.

For a clustering model we need to maximize the difference between the center of each cluster, so working on the variance, we need to use a pair of feature the less corelate to have a spread model.

Clusters are clearly appearing on graphs.