# AVILA

## Python For
## Data Analysis

Kevin Gomes

# DATA VISUALIZATION

# COMPOSITION

## Shape

20867 rows & 11 columns

## Target

Last column ('monk' → class)

## NaN

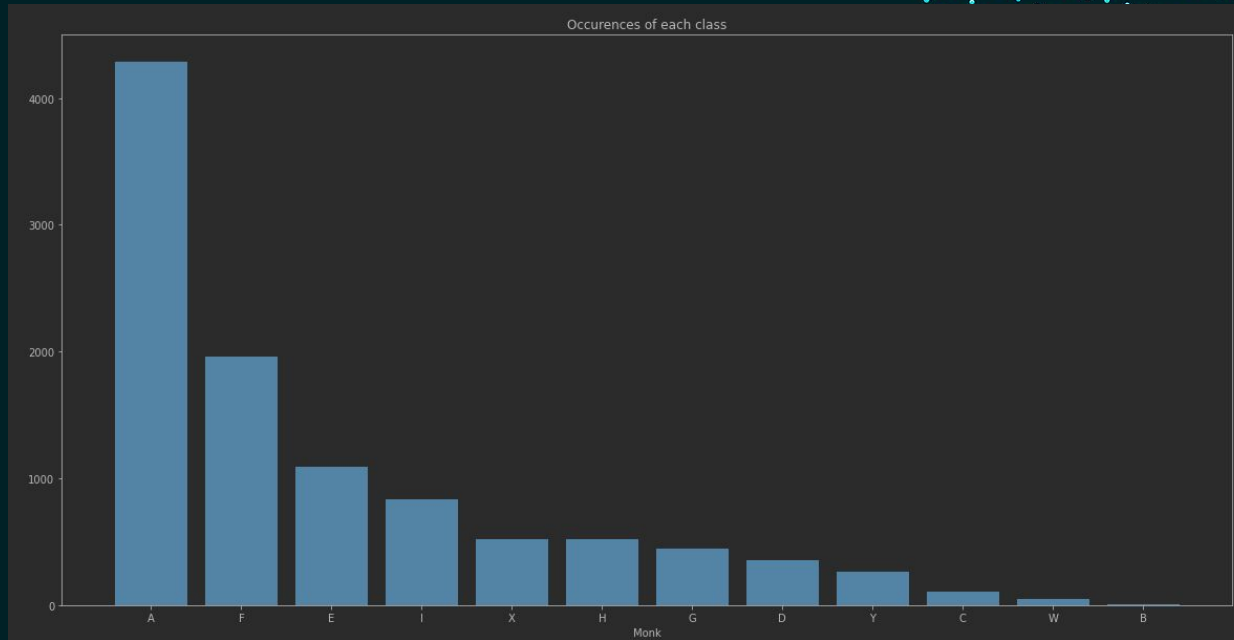No NaN in all the Data Set

# QUICK OVERVIEW

- Data already normalized by using the Z-Normalization method

- Divided in 2 data sets, train & test, with 10430 and 10437 samples respectively

Preview of the first 9 lines of the test set to see what the dataset looks like :

| | intercolumnar distance | upper margin | lower margin | exploitation | row number | modular ratio | interlinear spacing | weight | peak number | modular ratio/ interlinear spacing | monk |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.266074 | -0.165620 | 0.320980 | 0.483299 | 0.172340 | 0.273364 | 0.371178 | 0.929823 | 0.251173 | 0.159345 | A |
| 1 | 0.130292 | 0.870736 | -3.210528 | 0.062493 | 0.261718 | 1.436060 | 1.465940 | 0.636203 | 0.282354 | 0.515587 | A |
| 2 | -0.116585 | 0.069915 | 0.068476 | -0.783147 | 0.261718 | 0.439463 | -0.081827 | -0.888236 | -0.123005 | 0.582939 | A |
| 3 | 0.031541 | 0.297600 | -3.210528 | -0.583590 | -0.721442 | -0.307984 | 0.710932 | 1.051693 | 0.594169 | -0.533994 | A |
| 4 | 0.229043 | 0.807926 | -0.052442 | 0.082634 | 0.261718 | 0.148790 | 0.635431 | 0.051062 | 0.032902 | -0.086652 | F |
| 5 | 0.117948 | -0.220579 | -3.210528 | -1.623238 | 0.261718 | -0.349509 | 0.257927 | -0.385979 | -0.247731 | -0.331310 | A |
| 6 | 0.389513 | -0.220579 | -3.210528 | -2.624155 | 0.261718 | -0.764757 | 0.484429 | -0.597510 | -0.372457 | -0.810261 | A |
| 7 | 0.019197 | -0.040001 | 0.288973 | -0.042597 | 0.261718 | -1.013906 | 0.069175 | 0.890701 | 0.095265 | -0.842014 | F |
| 8 | 0.500607 | 0.140576 | 0.388552 | -0.637358 | 0.261718 | -0.681707 | 0.295677 | 0.931046 | 0.500624 | -0.642297 | H |

# CLASS DISTRIBUTION



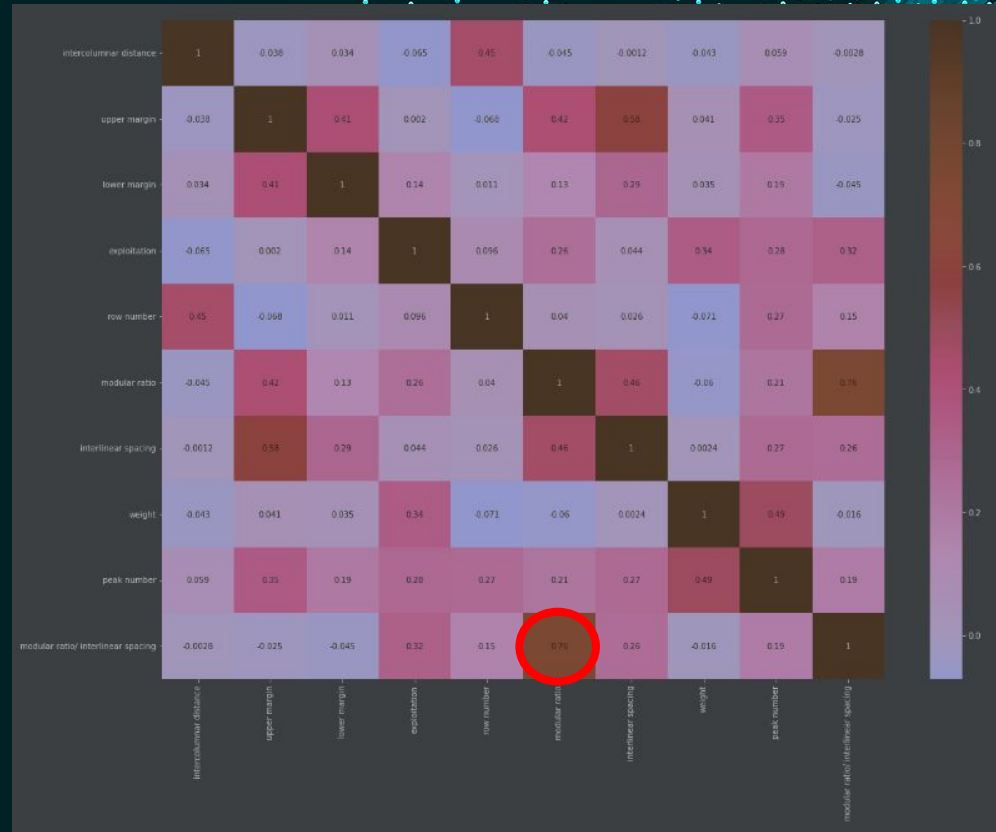The number of occurrences of each class shows an unequal distribution of classes.

While the most present class, A, appears more than 4000 times, class B appears only 5 times (training set).

# CORRELATION MATRIX

The Correlation Matrix is used to observe the links between the attributes.

In our case, the only 2 attributes that are too strongly correlated are 'modular ratio' and 'modular ratio / interlinear spacing' with 76% of correlation. (red circle)
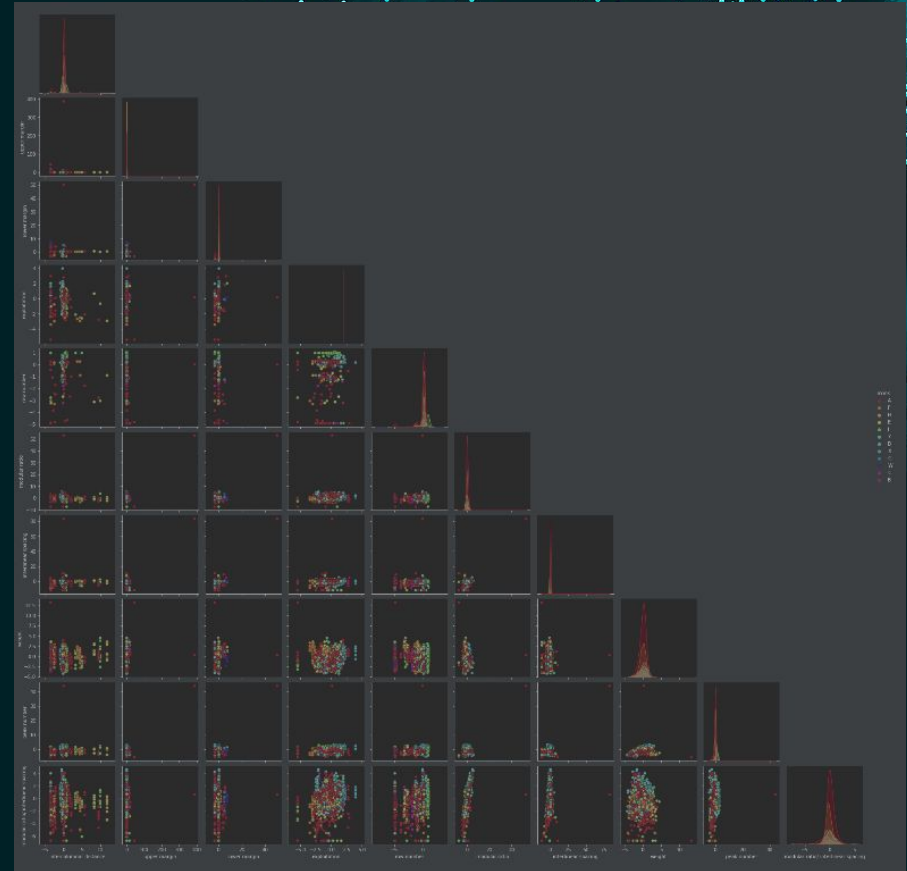
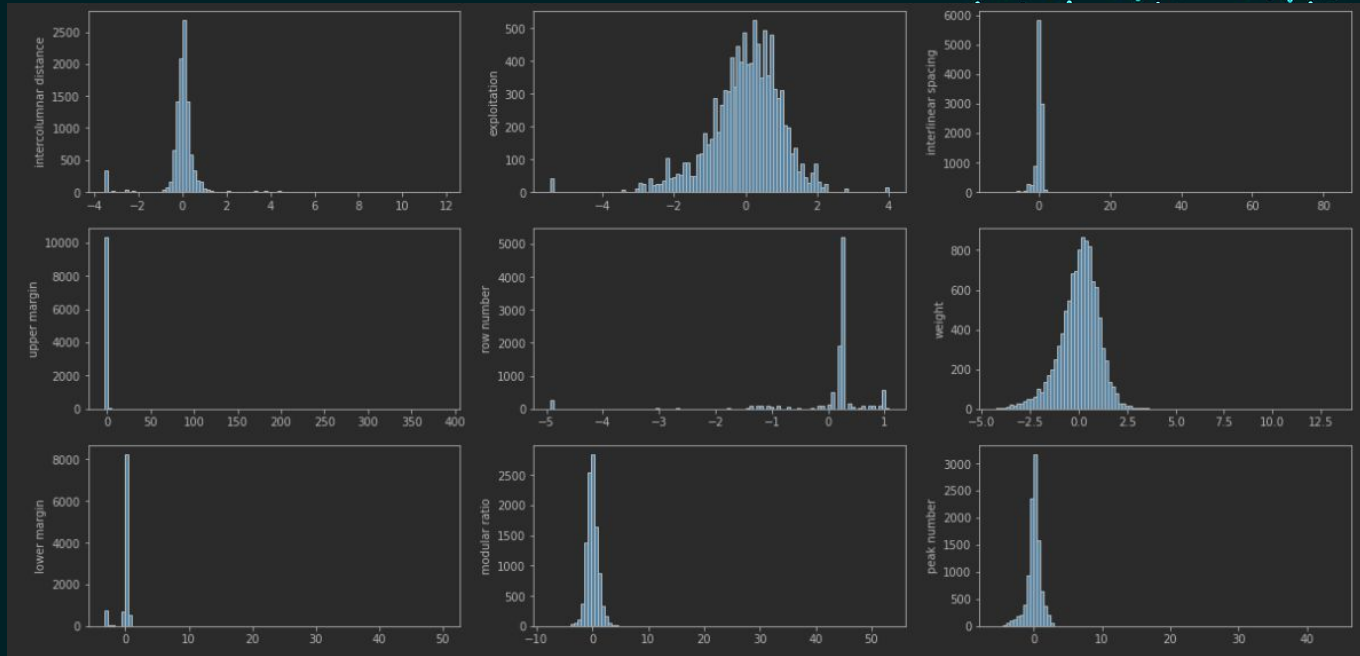To overcome this correlation, I decided to remove 'modular ratio / interlinear spacing' from the dataset.

# PAIR PLOT

Before deleting the attribute of the previous slide, I made pair plot between all attributes to observe in another way the correlation between attributes.

These pair plots confirm the correlation between our two attributes and therefore allow us to validate the deletion of the column 'modular ratio / interlinear spacing'.
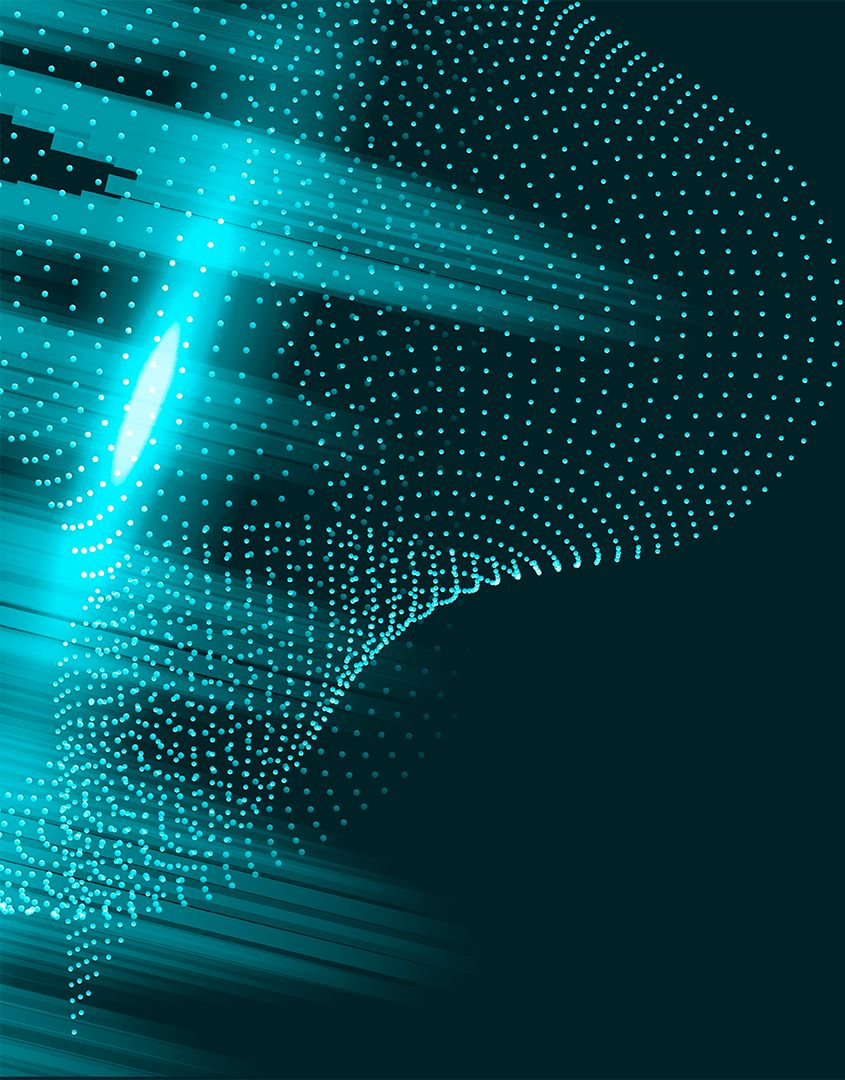
# ATTRIBUTES DISTRIBUTION



Thanks to the observation of the distribution of the attributes, we observe some extreme values for a majority of the attributes.
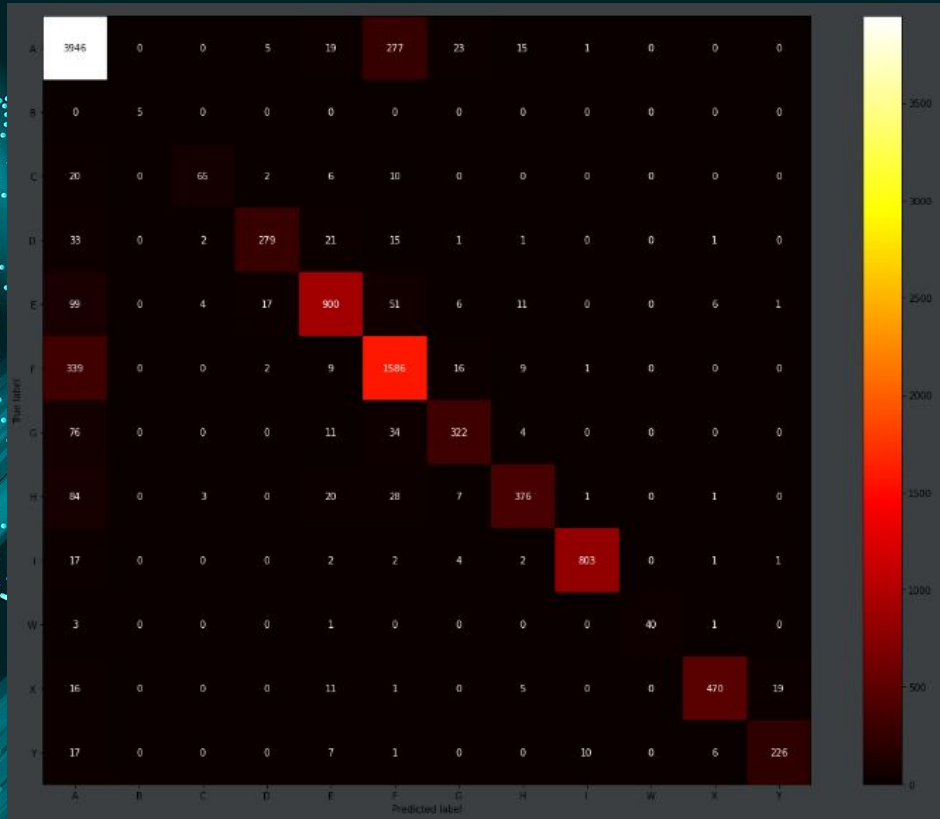
As for example 'upper margin' which has a graph that goes up to 400 whereas almost all the data are close to 0

MODELING

# K-NN



Confusion Matrix

Best Score : 0.837967

Works better with :
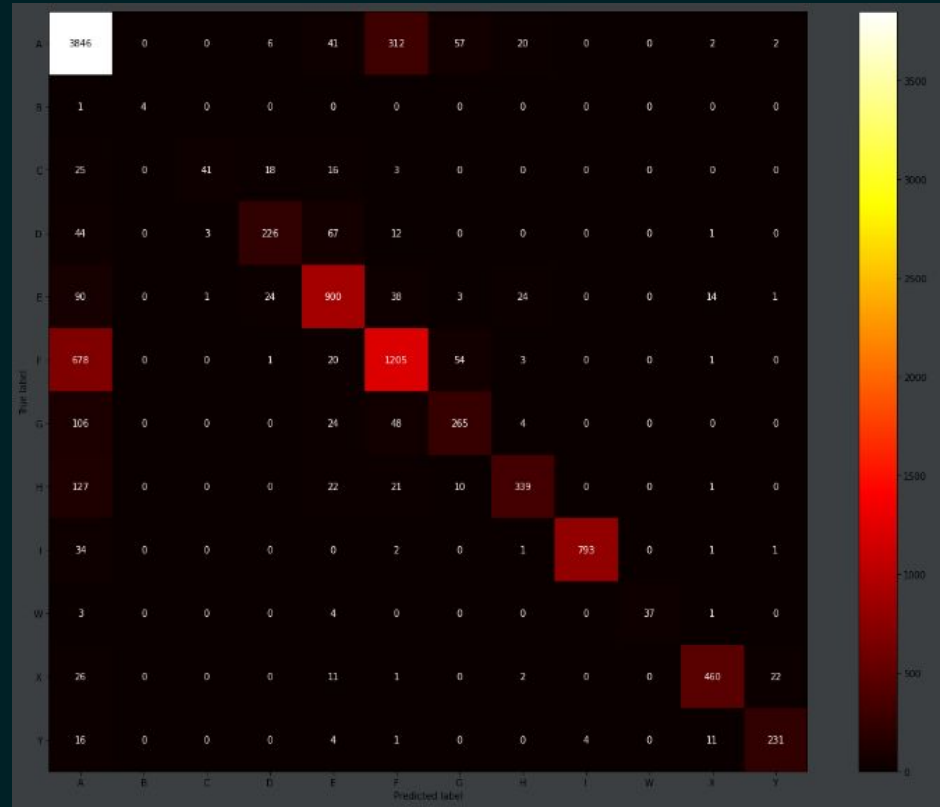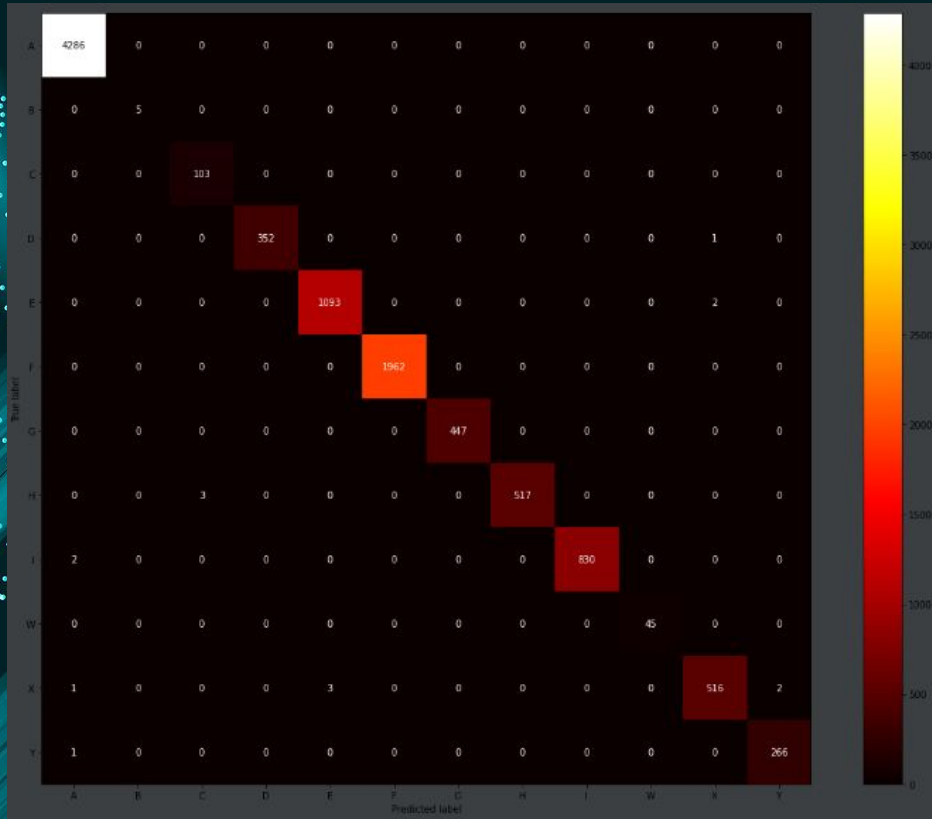- leaf_size = 15
- n_neighbors = 5
- p = 1

# SVC

Best Score : 0.788591

Works better with :
- C = 2
- gamma = 0.5
- kernel = rbf



Confusion Matrix

# XGBOOST


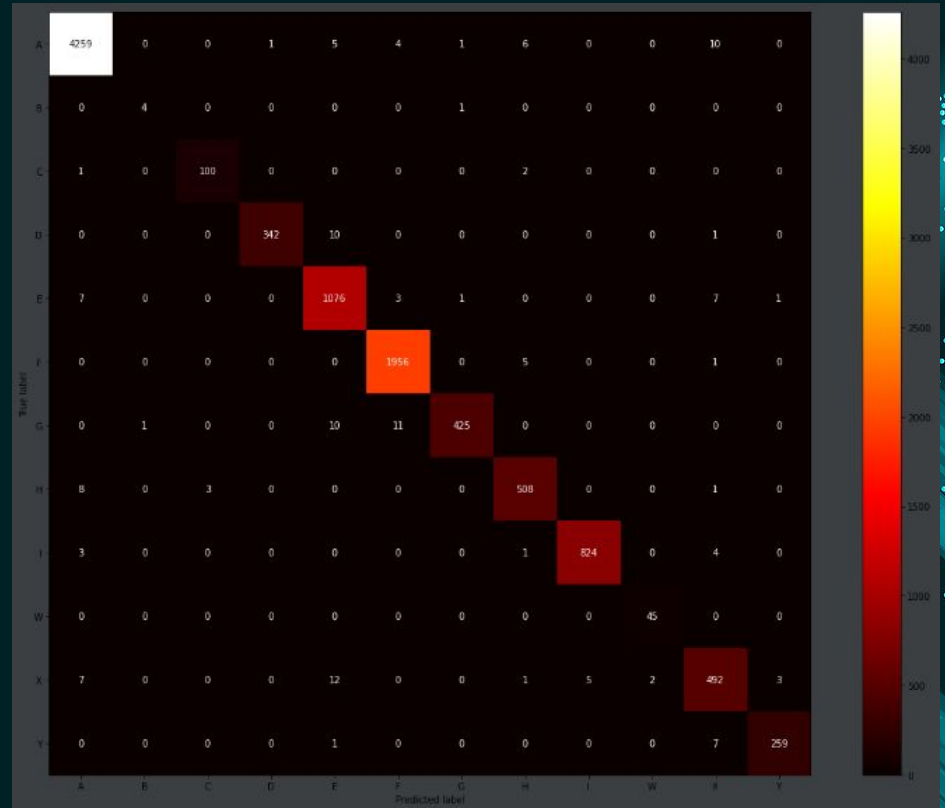Confusion Matrix

Best Score : 0.996357

Works better with :
- learning_rate = 0.05
- n_estimators = 1500

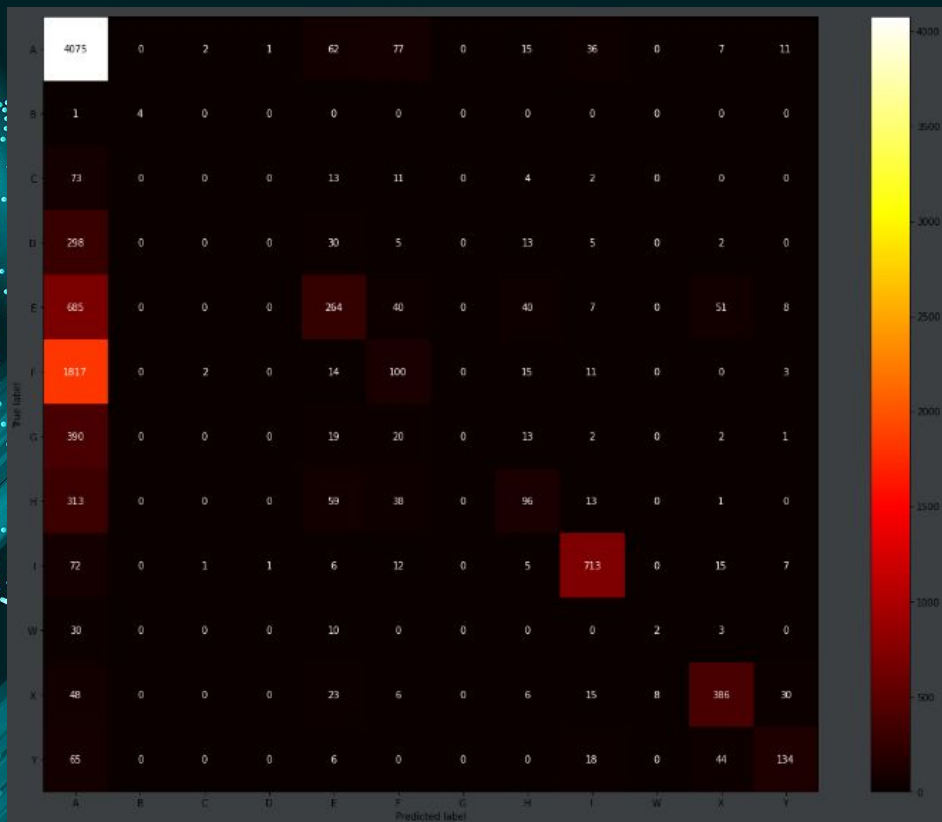# DECISION TREE

Best Score : 0.982933

Works better with :
- criterion = entropy
- max_features = 9



Confusion Matrix

# LOGISTIC REGRESSION



Confusion Matrix

Best Score : 0.552733

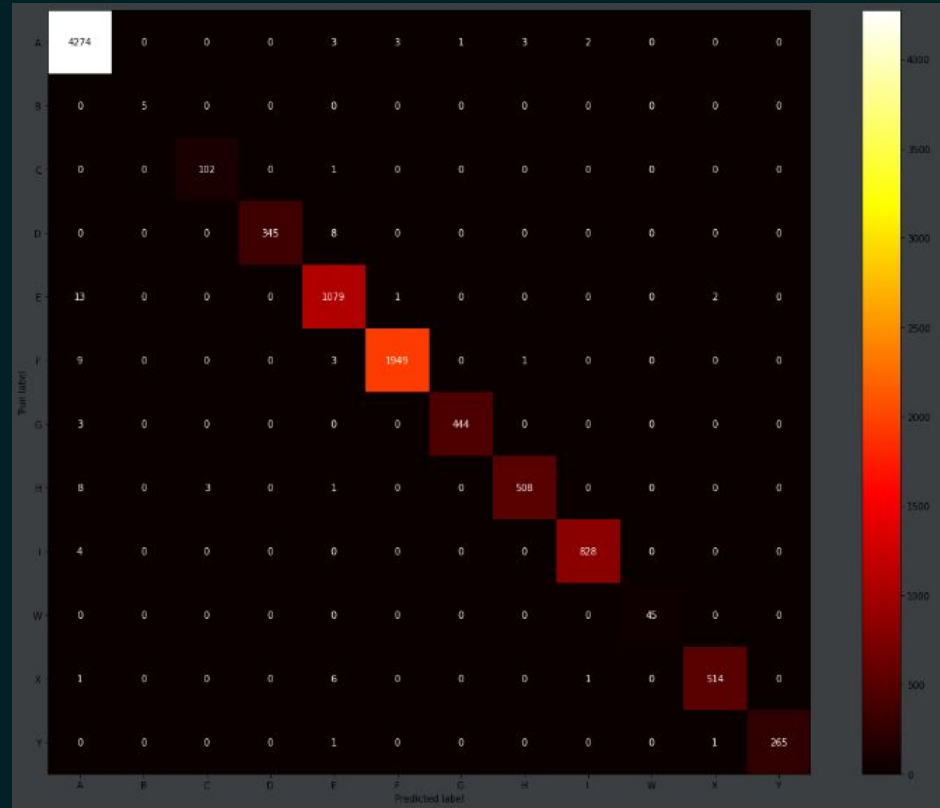Works better with :
- max_iter = 500
- nb_jobs = 2

Confusion Matrix

# COMPARISON

| Model | Score |
|---|---|
| XGB Classifier | 0.996357 |
| Random Forest | 0.988878 |
| Decision Tree | 0.982934 |
| KNN | 0.837967 |
| SVC | 0.788591 |
| Logistic Regression | 0.552733 |

- XGB Classifier seems to be the best model but Random Forest & Decision Tree algorithms are also very accurate.

- KNN & SVC still have pretty good results.

- Logistic Regression has an accuracy of a little over 50% which is not very good.

# CONCLUSION

# CONCLUSION

This project allowed me to better understand how to work on a data set, how to exploit it and predict its data by comparing several models.

This project, as a whole, went well and I am satisfied with the results I obtained during my predictions.

My only regret is that I didn't have the time to make an API for my project due to health problems (COVID).