# ESILV - Python for data analysis - project 2022

Thomas BAUDU and Paul CANAL

# Introduction

This year our project in Python for data analysis was to analyse a dataset and create a model in Machine Learning.

To analyse the dataset, we used several graphs like barplots, boxplots...

For the part in machine learning, we tried 4 different models that we will present you later.

To present our finished work we made an API using Flask, you can find more info on how to use it on our [GitHub](#)

# Summary

Introduction


The Avila bible

Analyse

Model in Machine Learning

Flask API

Conclusion

# The Avila bible

Italian Romanesque manuscript whose text and images were completed after his arrival in the Cathedral of Avila at the end of the 12th century.
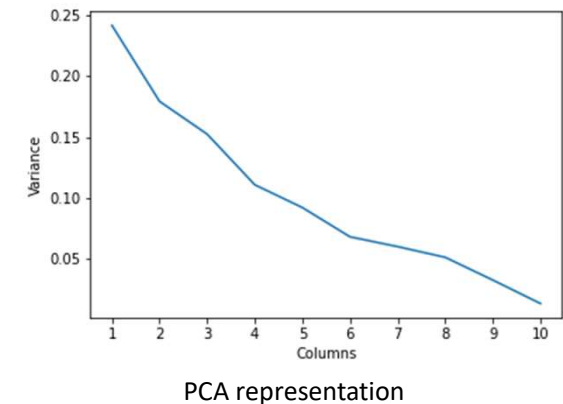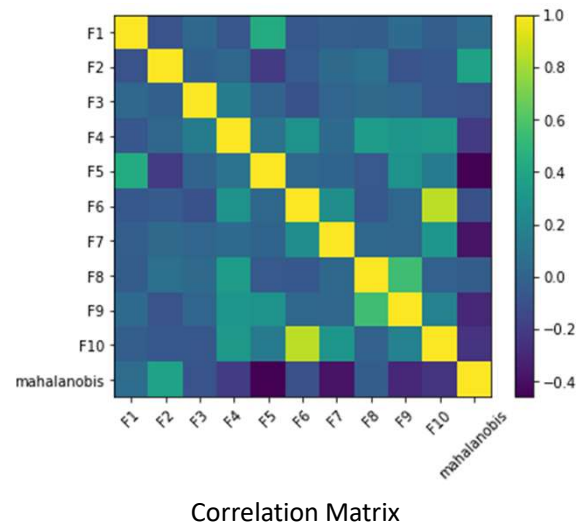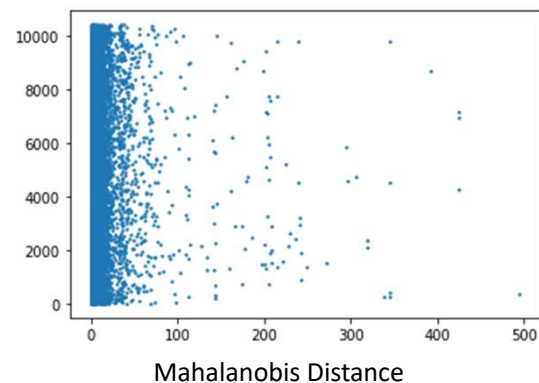
The Avila Bible is a copy of the Latin Vulgate which, due to its large dimensions -842 pages of 580 x 390 mm in parchment and about 15Kg in weight- is considered the best exponent in the group of "Atlantic Bibles"large manuscripts for the liturgical use of monastic communities, created between the end of the 11th century and the end of the 12th century, in the midst of a Gregorian reform and in the geometric style very different from other European models, which was already common in Roman mural painting in the 12th century, and which is characterized in the miniature by the initials that start each Biblical book, with geometric frets and a very sophisticated decoration, which we also find in other Italian codices of that period, especially at the desks in Rome and Milan.

Although there is no record of its initial dating or origin, there is no doubt that its origin would be Italian from the middle of the twelfth century and that it would reach Avila around 1175 where the work was continued, in a style very different from that of the Italian part.

# Analysing the dataset

The first graph represents the Mahalanobis Distance of the DataFrame. There is 10428 point on this graph. But only a hundred have a Mahalanobis distance > than 100. So, we decided to delete all the points with a distance > than 260 to delete the inconsistent points.

The second graph is a Correlation Matrix in the dataset. We can observe which columns are single and who help us to guess on the test set. The last one have the same goal, this is the PCA Method.

Mahalanobis Distance

Correlation Matrix

PCA representation

# Model in Machine Learning

For our model in Machine Learning, we tried 5 different models.

The five graphs below are the confusion matrix of each model. We easily see that the two logistic regression's confusion matrix are not as good as the others. The best ones are from the Random Forest models, and the one without the PCA gives the best results.



Logistic regression



Logistic regression PCA



K-Nearest Neighbors



Random Forest



Random Forest PCA

# Model in Machine Learning

Finally, here are some important rates for each model :

| Model | Accuracy | Precision | Recall | F1_Score |
|---|---|---|---|---|
| log_reg | 0.561751 | 0.481157 | 0.4025752 | 0.4089718 |
| log_reg_pca | 0.562135 | 0.480602 | 0.4028264 | 0.4090753 |
| knn | 0.757497 | 0.812033 | 0.6929663 | 0.7413758 |
| rand_for | 0.985245 | 0.990682 | 0.9794742 | 0.9849854 |
| rand_for_pca | 0.77963 | 0.865399 | 0.6655743 | 0.7324215 |

We clearly see that the Random Forest model (without the PCA) gives the best results for every rate

PCA seems to not have a lot of influence for the logistic regression model, however it is the case for the Random Forest model and we did not except the model without the PCA to be better.

# Flask API

To represent our model, we made an API using Flask. It consists of a Python application and some HTML pages, along with CSS files.
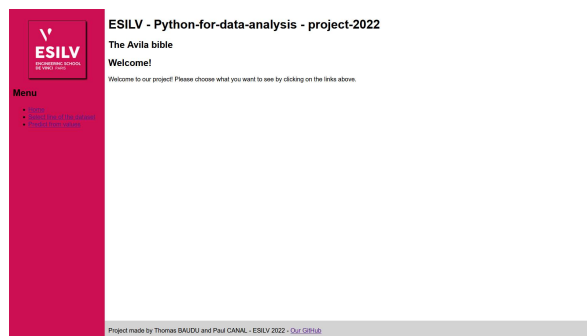
In the Python application, we define functions to be called when we access a web page.
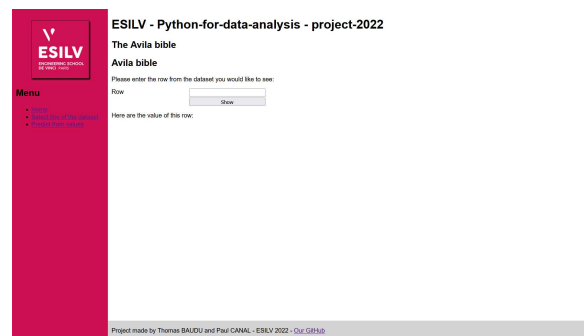The first page is index.html, where we just present the project.
On the /avila_bible path (show_bible.html), we made a function to show the specific row of the dataset the user asked to see.
On the /predictions path (predict.html), we made a function to use the model we made: the user enters specific values and we predict which copyist would have had those values.
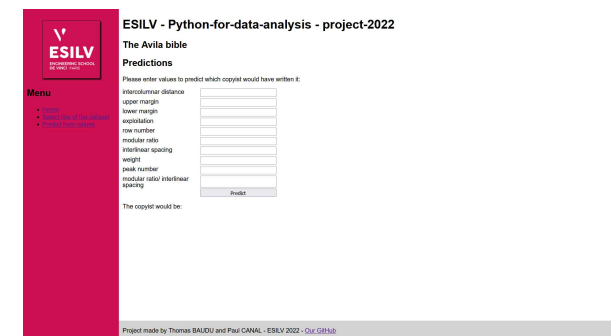
Each HTML page is an addition to the base.html file, to keep some coherence between all pages.



*index.html*



*show_bible.html*



*predict.html*

# Conclusion

The dataset we choose was interesting to analyse and also original compared to the other datasets (like detecting spam or drug consumption).

With around 0.985 of accuracy, we have a really good model. This result was achieved quite easily.

Flask was new to both of us so it was a new thing to learn, we also never used API so it was a good thing to start. It is also not too complex and we can made a lot of things with it.