# Python Project

UNDERWOOD Theo and TOUBIANA Benjamin

DIA 7

2021-2022

# Problem : Determine if a chemical is ready biodegradable

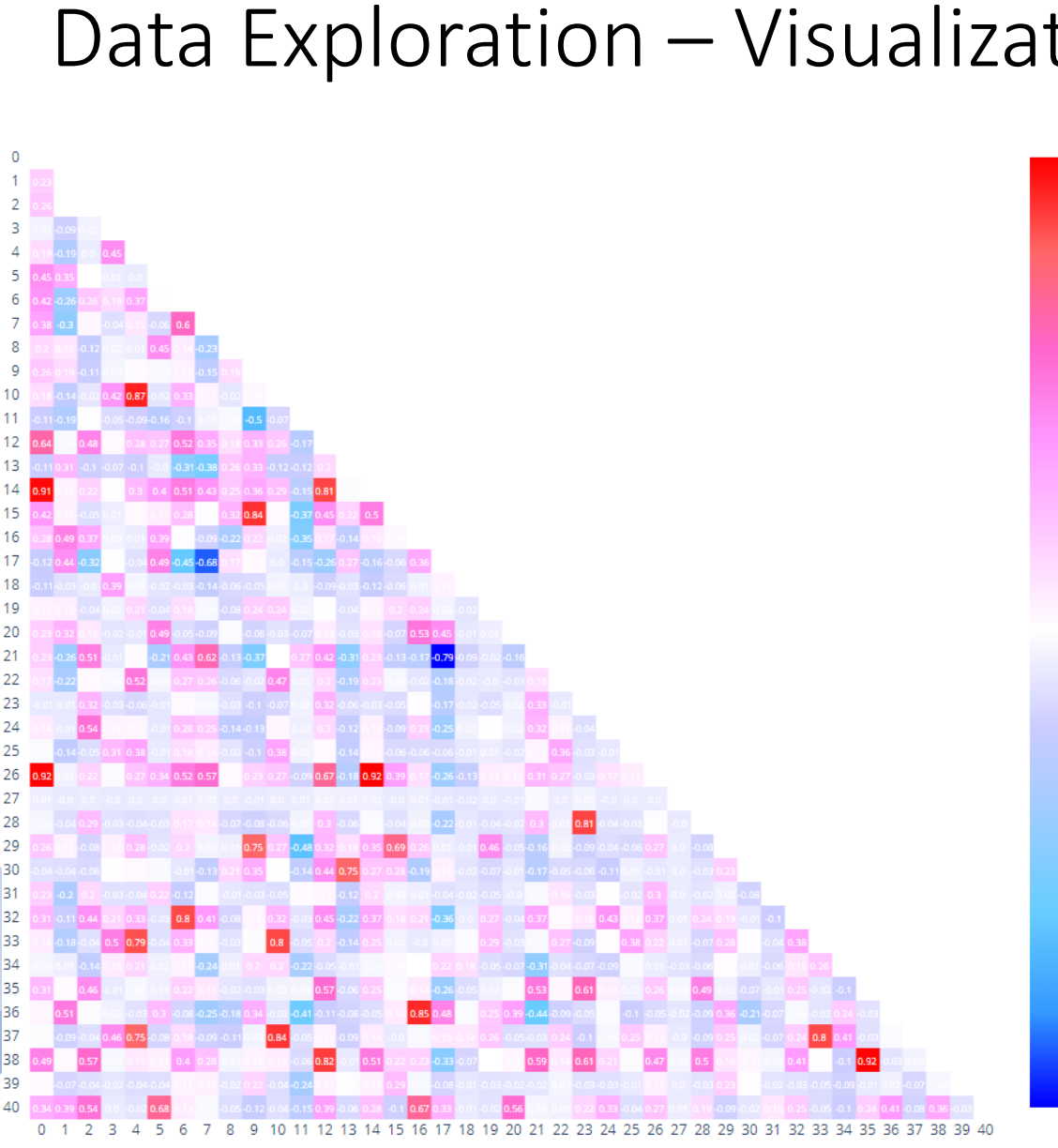What : **QSAR biodegradation Data Set**

How: Molecular characteristics of each chemical

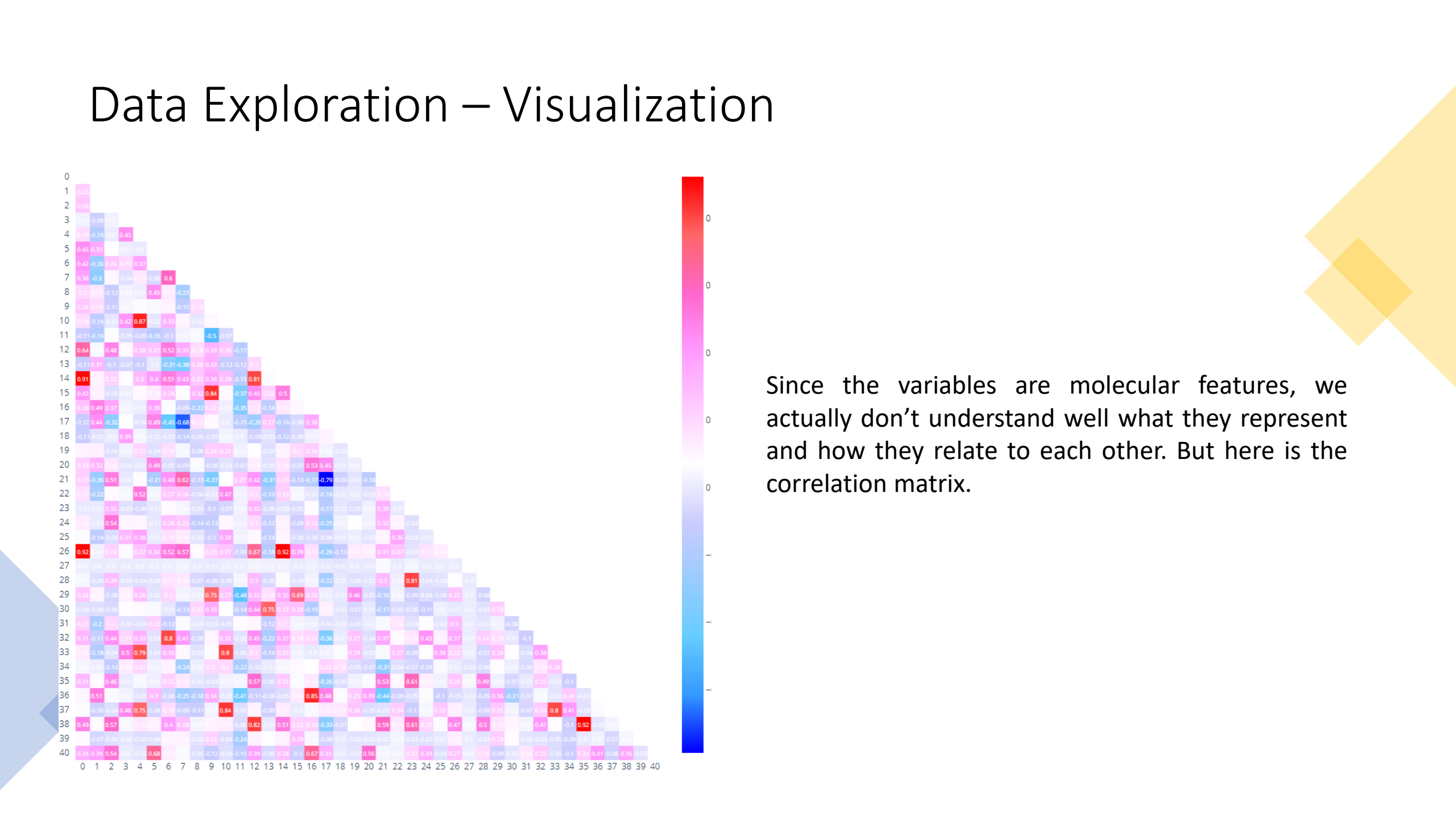Why : Classify chemicals by ready biodregradability

# Data Exploration – Raw Data

- File format : ''.csv''

- Individuals : 1055 Chemicals

- Features : 41 Molecular characteristics

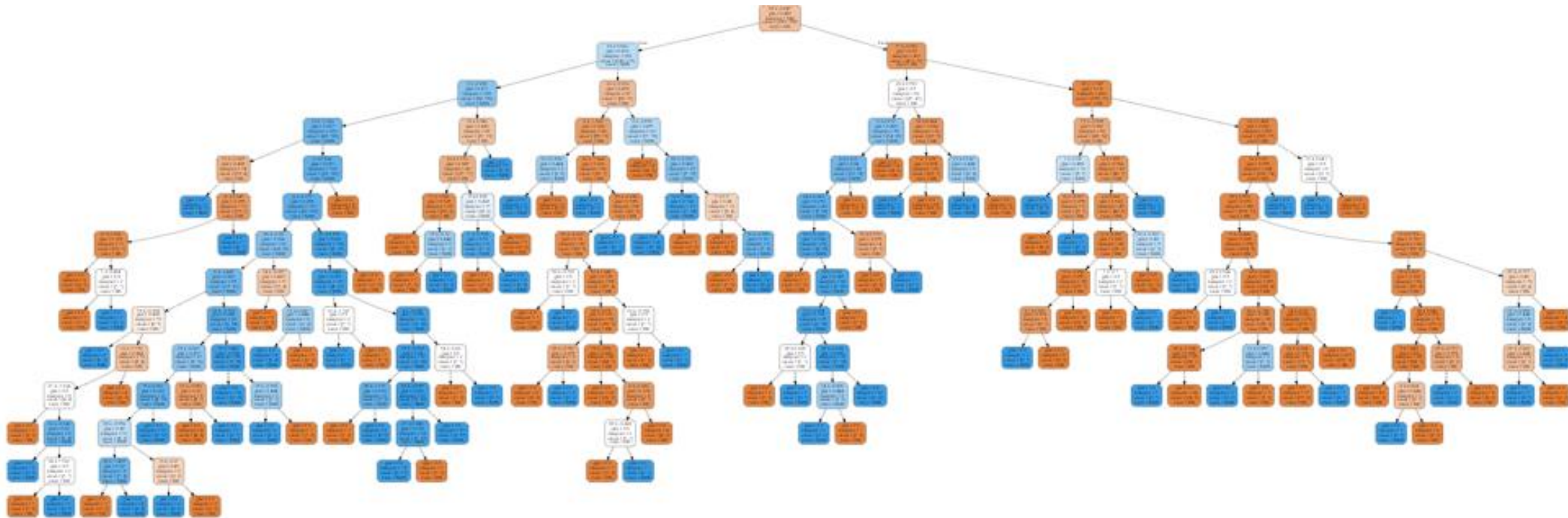- Targets : 2 possibilities – Ready Biodegradable, Non Ready Biodegradable

# Data Exploration – Visualization



Since the variables are molecular features, we actually don't understand well what they represent and how they relate to each other. But here is the correlation matrix.
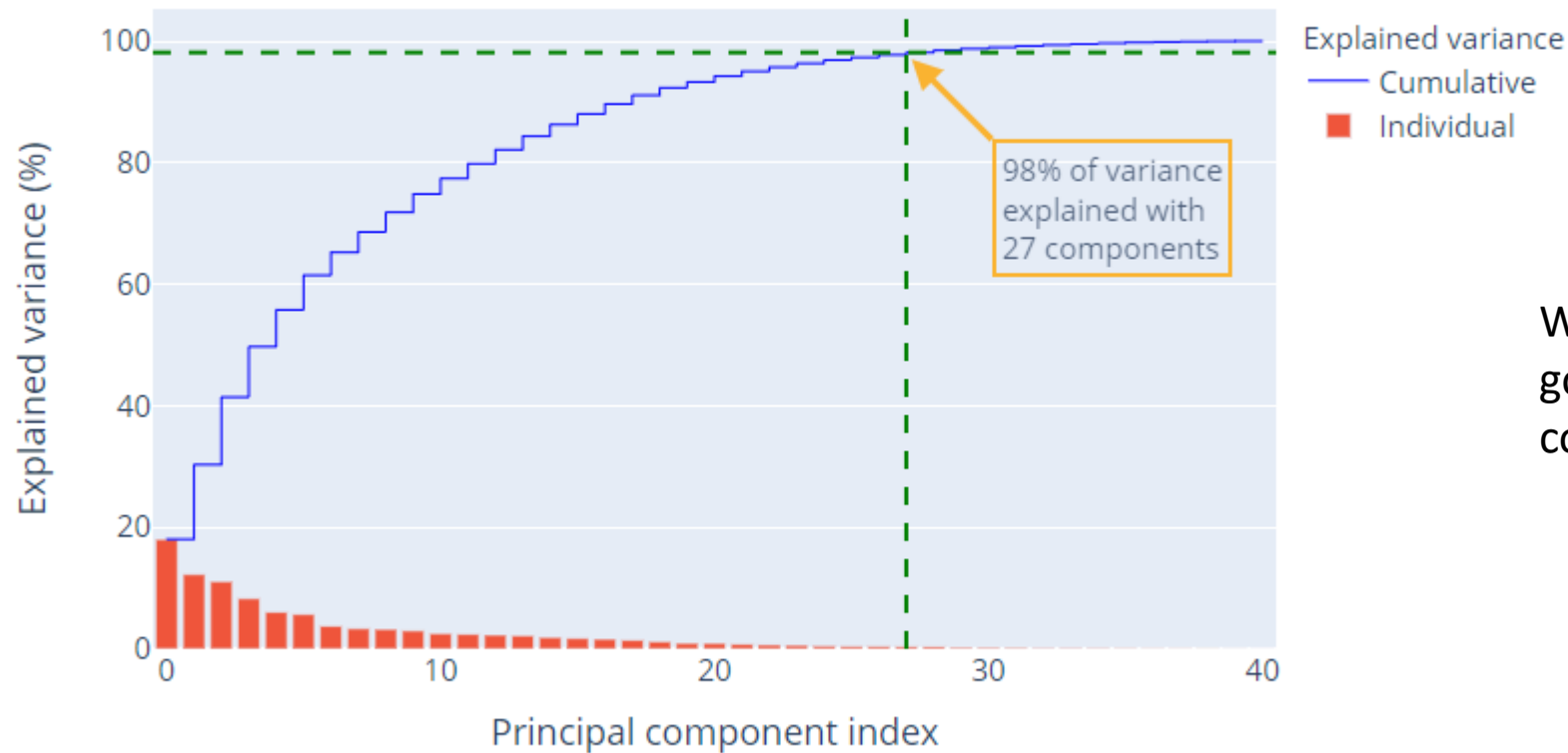
- Here is a plot of a decision tree to find out which features have the greatest impact on the predictions. It's not very clear because of the many variables.

]:

# Data Pre-processing

- We did not change the base data since there isn't missing values

- Standardization of the data (mean of 0 and standard deviation of 1)

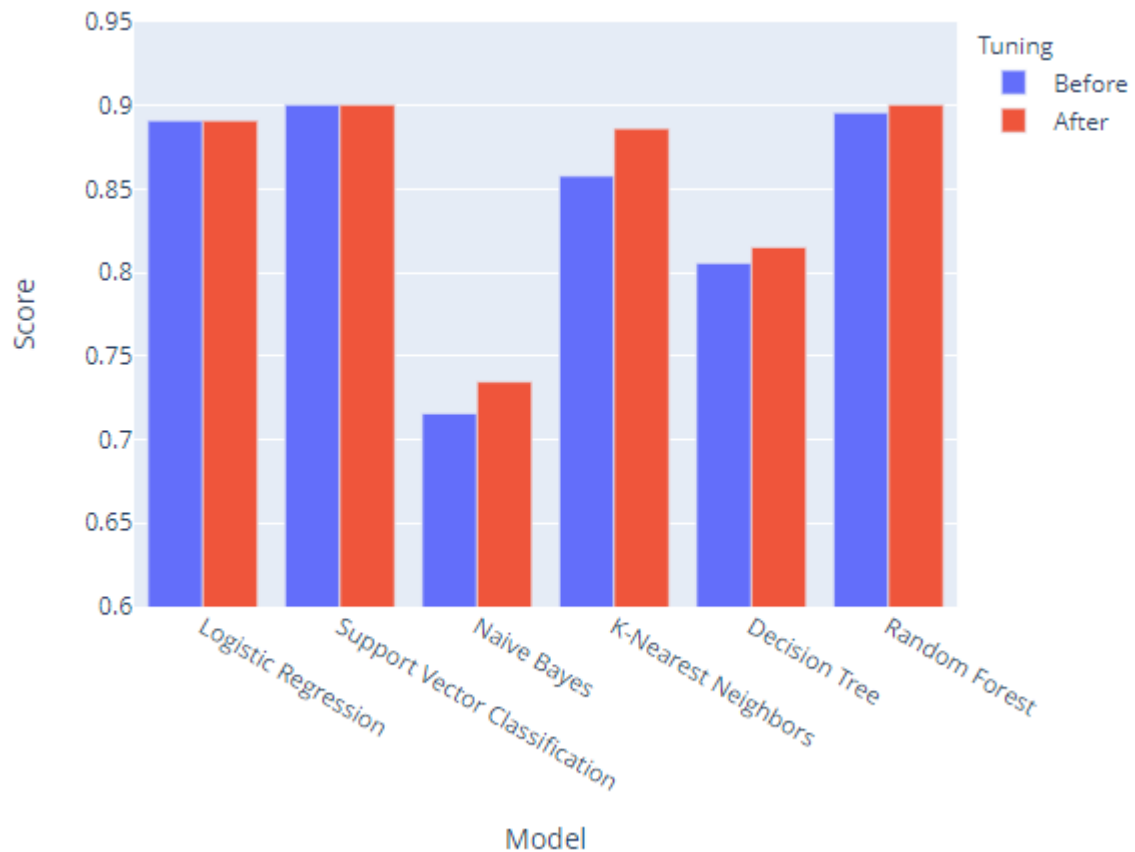- Splitting dataset into training and testing sets

We made an ACP to find how to go about reducing the amount components we need.

We obtain a variance of 98% with 25 components, so we finally decided to keep all the variables because treating 27 or 41 is the same for us. We would have used PCA if we had obtained less than 10 components.

# Modelization

| Model | Accuracy before tuning | Accuracy after tuning |
|---|---|---|
| Logistic Regression | 0.89 | 0.89 |
| Support Vector Classification | 0.90 | 0.90 |
| Naive Bayes | 0.72 | 0.73 |
| K-Nearest Neighbors | 0.86 | 0.89 |
| Decision Tree | 0.81 | 0.82 |
| Random Forest | 0.89 | 0.90 |

# Conclusion



- The hyperparameters tuning has often improved the accuracy of the models even if sometimes it was the same.

- We notice that the ranking before and after the tuning remains the same, which confirms that some models are better than others in some situations.

• The type of data processed is numbers and the site of scikit-learn recommends SVC for this kind of dataset. That's why it's the best model.



scikit-learn algorithm cheat-sheet