



# Python for Data Analysis

Jules THUILLIER – Victor THUILLIÉ  
DIA5

# Context

This dataset summarizes a heterogeneous set of features about articles published by Mashable in a period of two years.

The goal is to predict the number of shares in social networks (popularity).

There are 61 variables, 58 predictive, 2 non predictive and the target variable

# Context

As a first observation, it would seem that there are only numerical values

At first, we will try to predict a continuous value (the number of shares), so we place ourselves in a case of regression

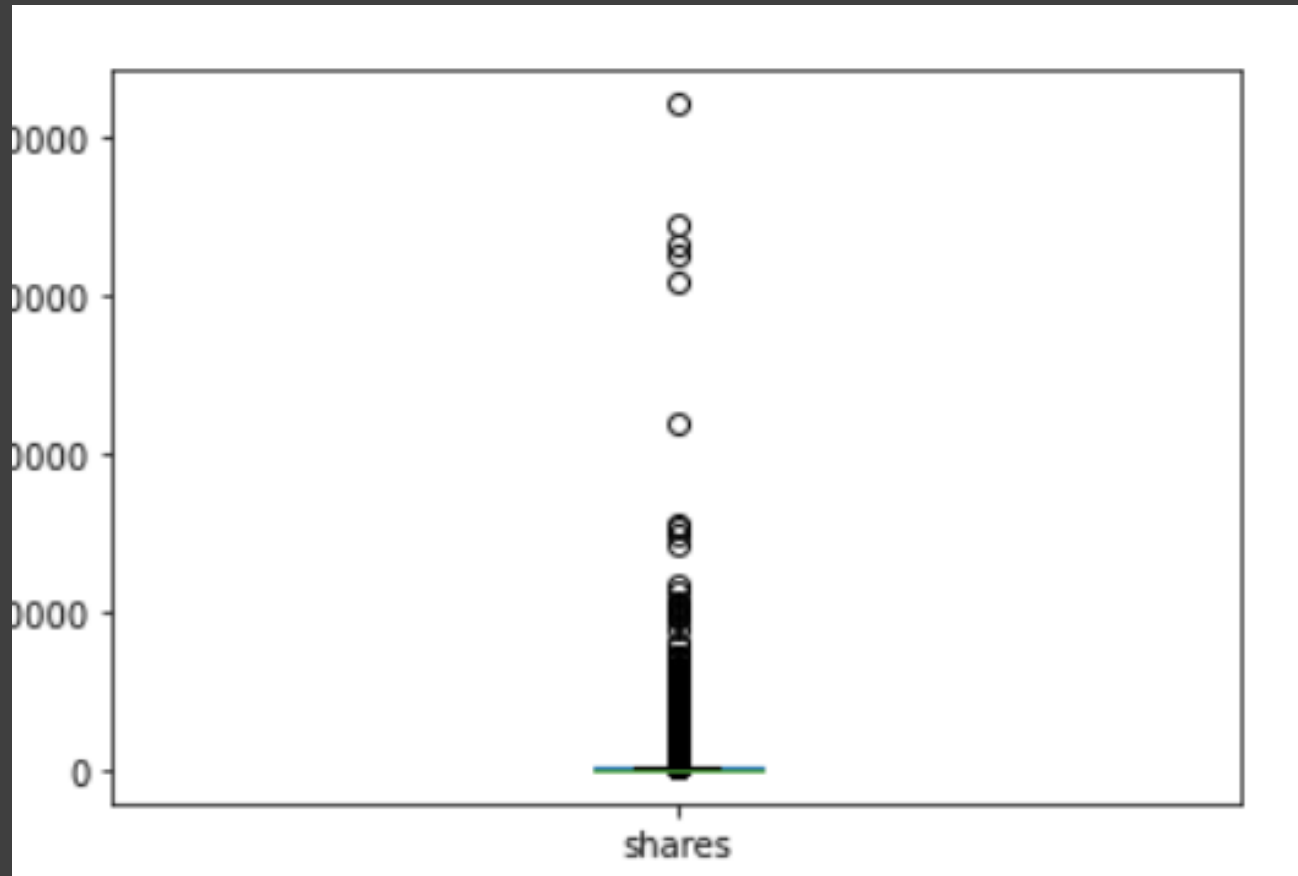
However, we can already add the 'popular' column for later. ( We will explain it later )



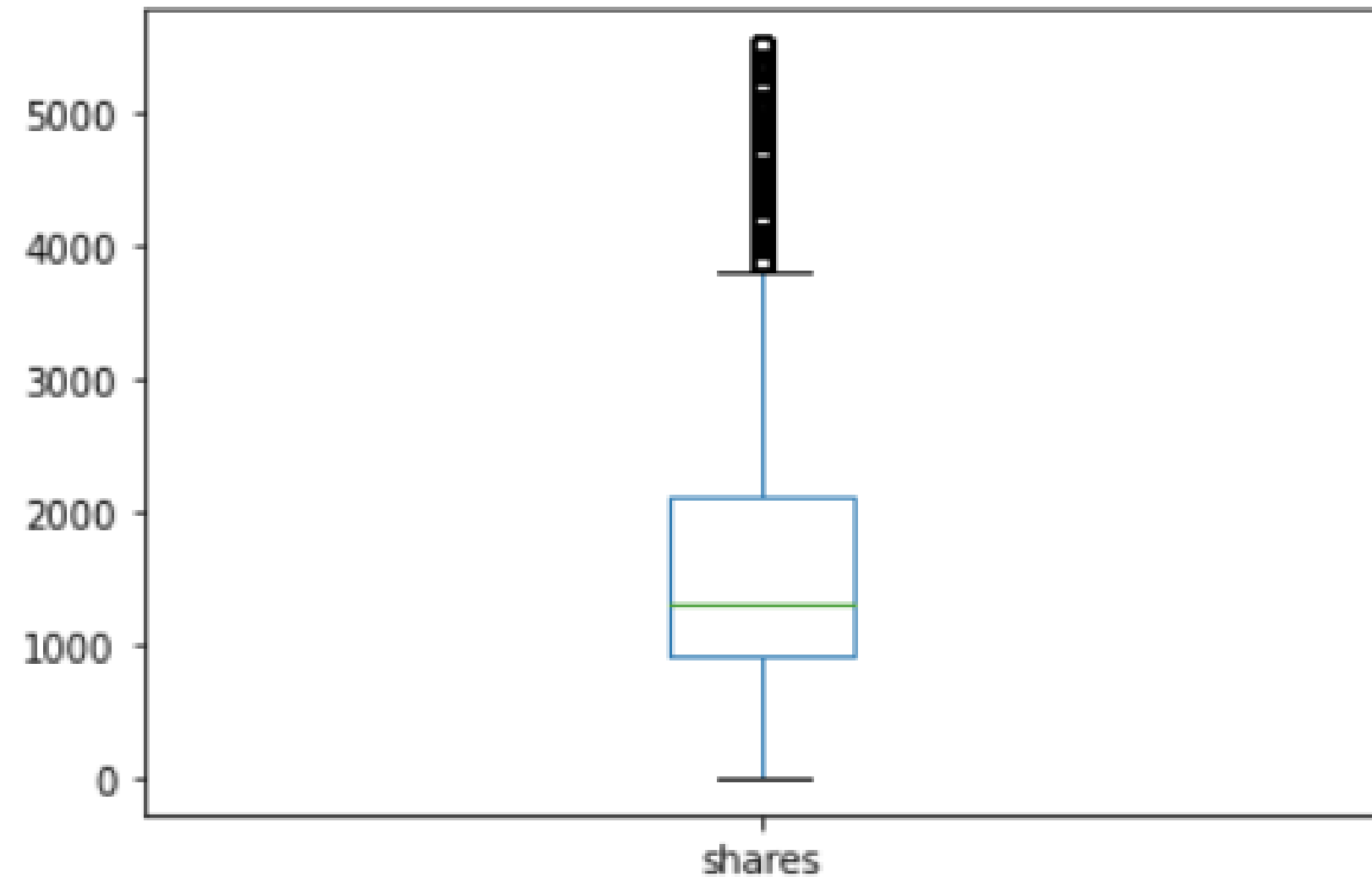
# Table of Content

- [Context](#)
- [Data Exploration](#)
  - [Outliers values](#)
  - [Correlations between values](#)
  - [Missing values](#)
- [Regression](#)
  - [Feature selection](#)
  - [Data Modeling](#)
- [Classification](#)
  - [Feature Selection](#)
  - [Data Modeling](#)
  - [Hyperparameter Tuning](#)
    - [GridSearch CV](#)
    - [Learning Curve](#)
  - [Model Export](#)

# Data Exploration – Outliers Values



- We can see that the boxplot is clearly flatten by the outliers values
- That's why we clean the data by removing it



## Data Exploration – Outliers Values

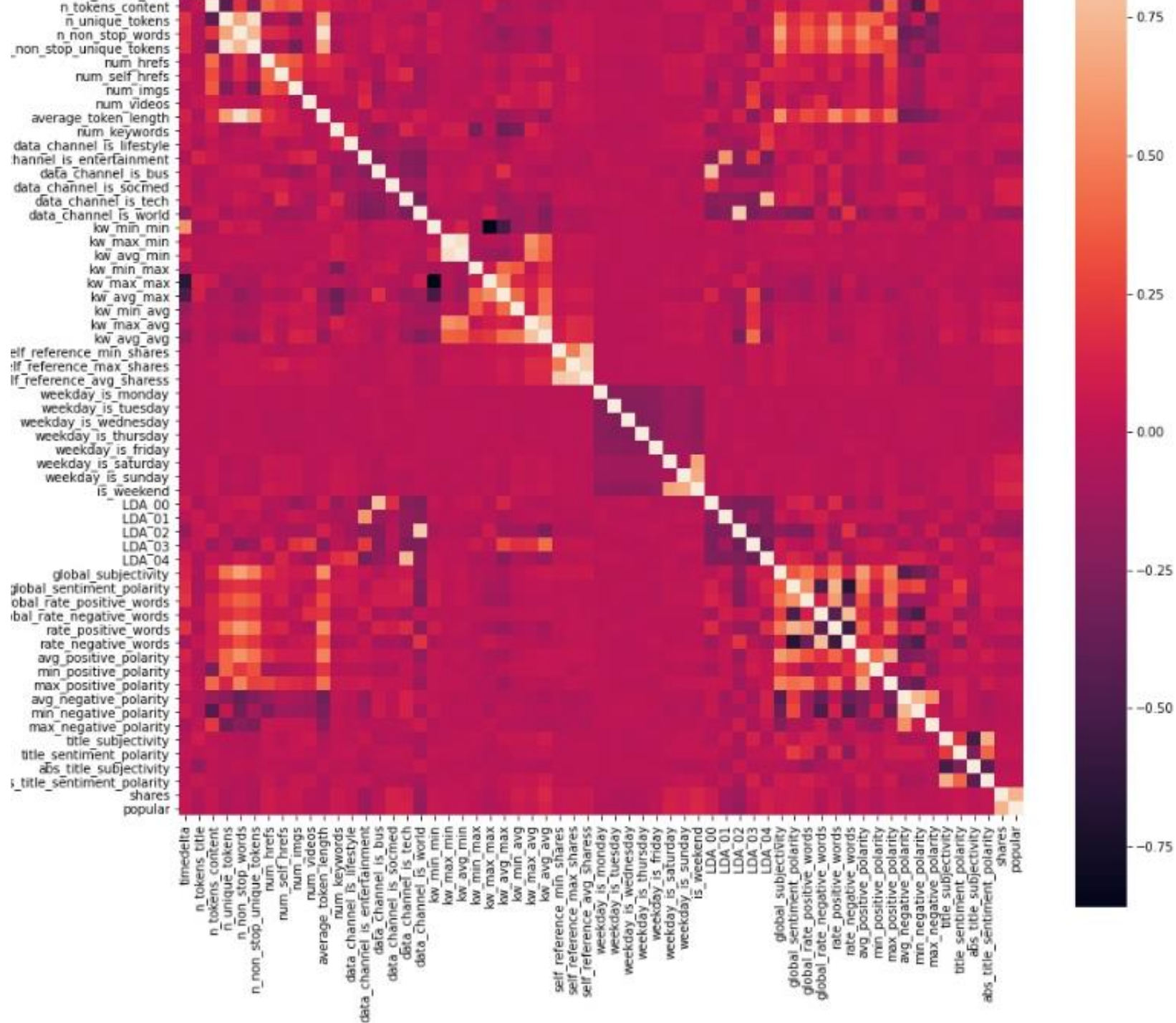
Data state after removing outliers

# Data Exploration – Correlations Between Values

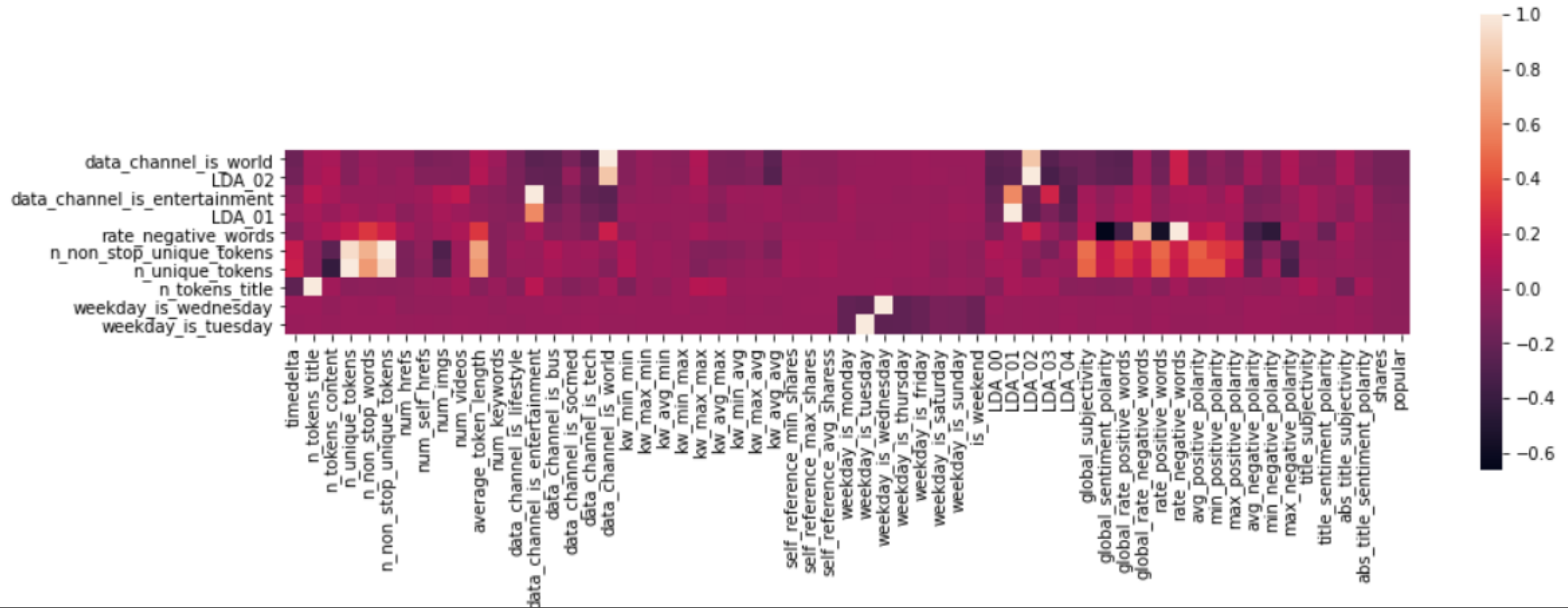




- There are a lot of variables but if you look closely, none of them seem to be correlated with the target variable.
- As this heatmap is not easy to read, we have two more, highlighting the 10 variables most correlated with the target variable and the 10 least correlated.

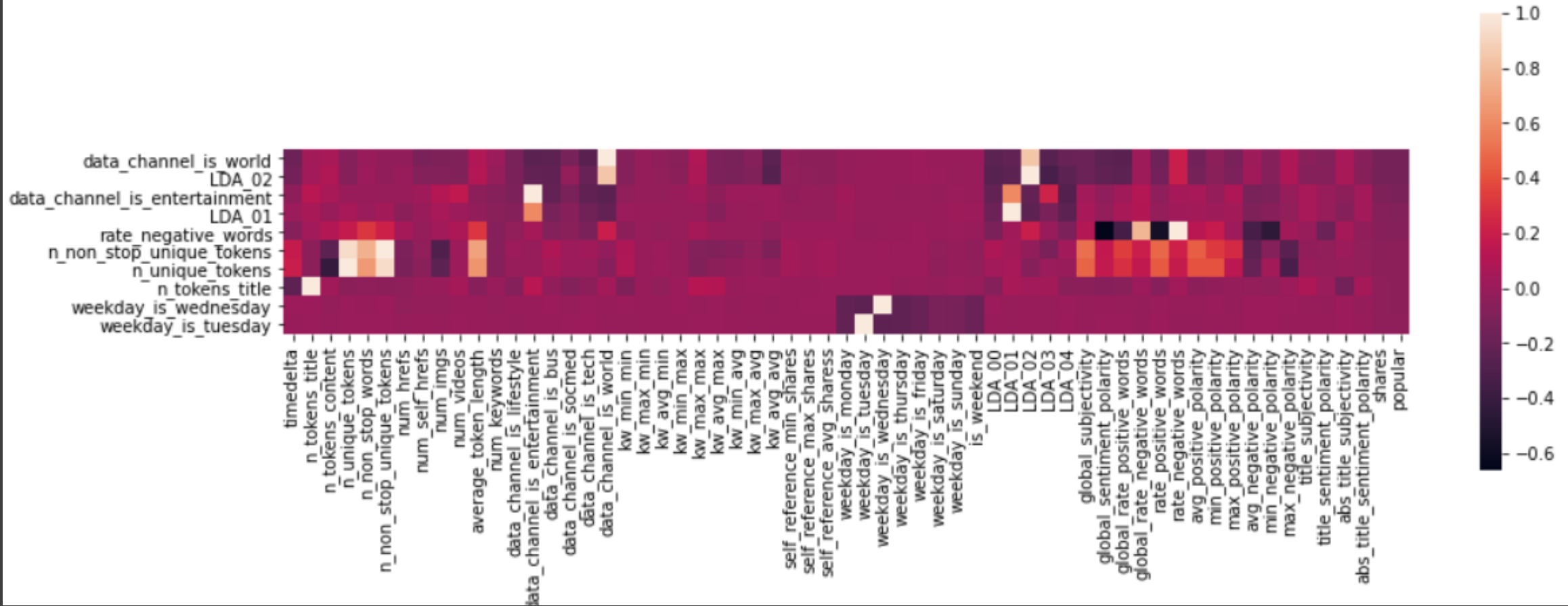






## Data Exploration – Correlations Between Values

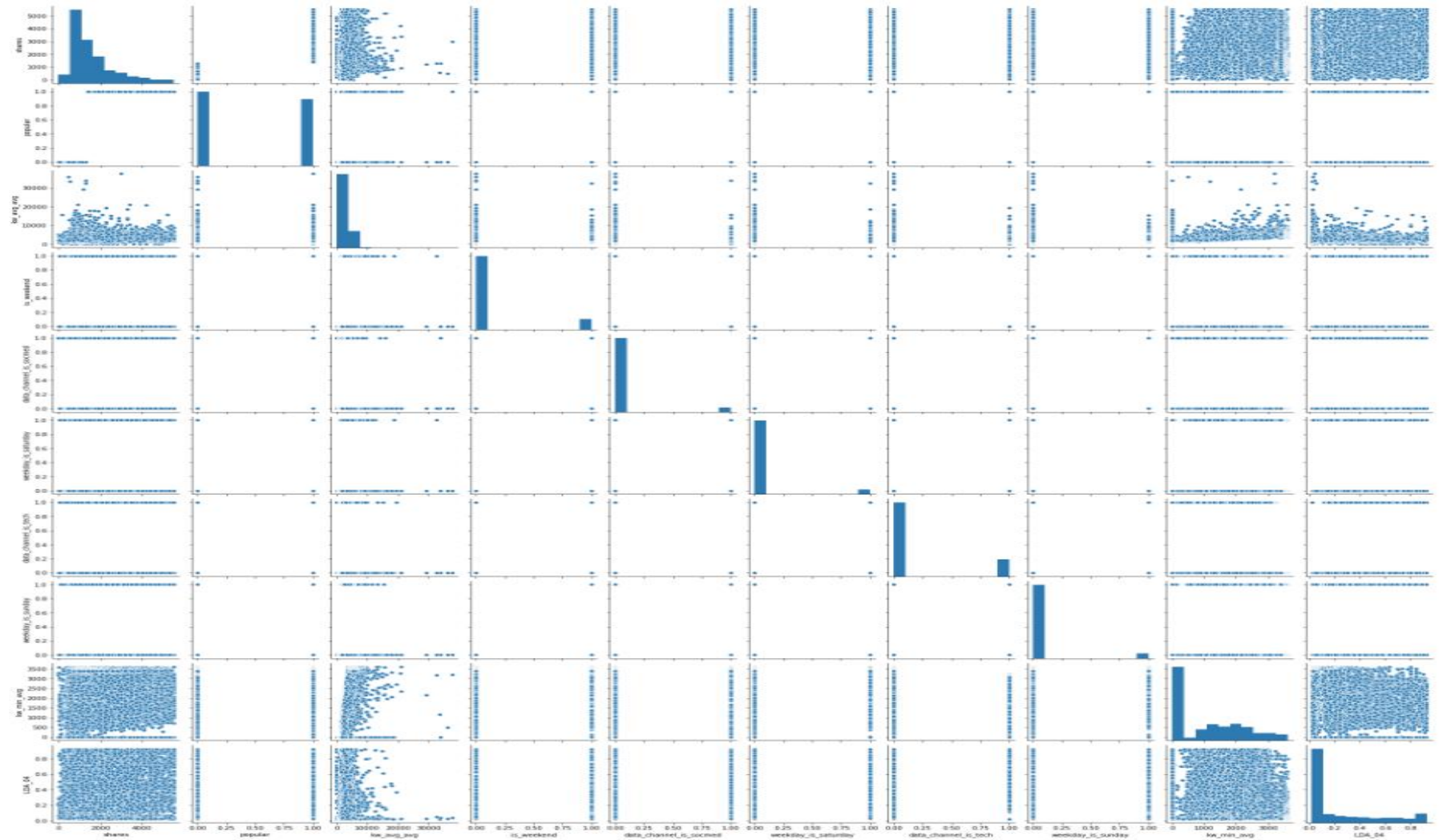
- On this heatmap we see that the 10 variables with the lowest correlation with the target variable do not have a significantly negative correlation.
- The lowest correlation value is around -0,14



## Data Exploration – Correlations Between Values

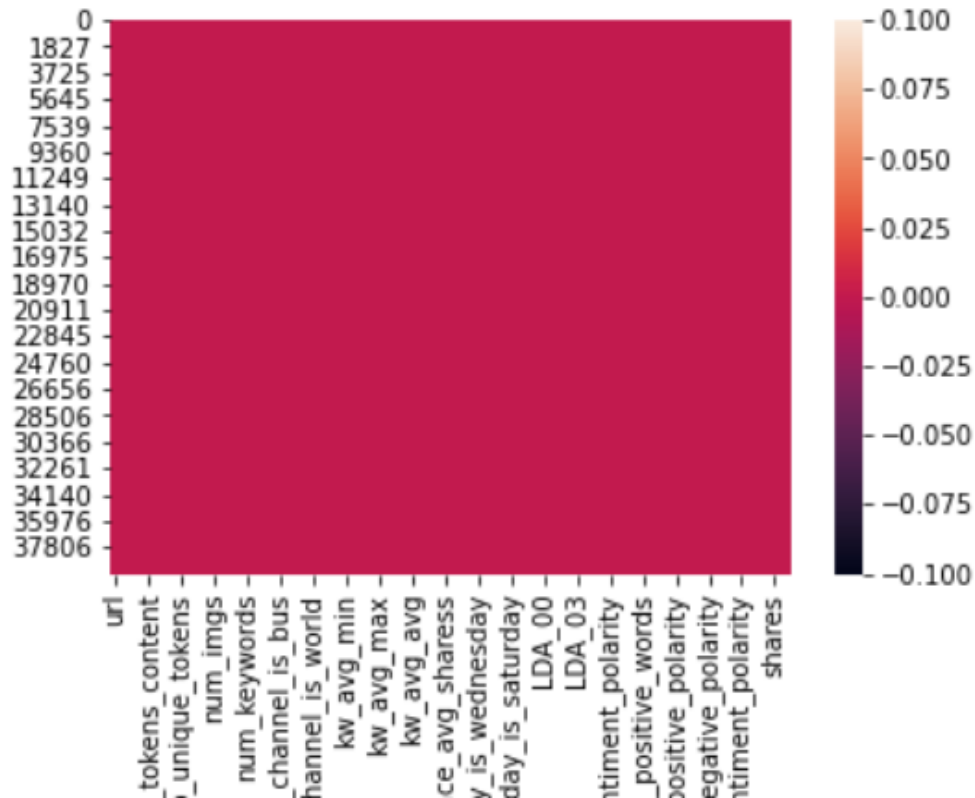
- Here we have the 10 variables most correlated with the target variable, as for the previous heatmap, we find that none of these variables are really correlated with the target variable.
- The highest correlation value is around 0,04

- Shares are on the first line, even if it's quite difficult to read, we can clearly see that there is no really pattern or interesting pattern between variables that have the highest correlation with the target variable.
- We only have cloud of points or just a separated values.





# Data Exploration – Missing Values



- By looking at the description of the dataset and this heatmap we can assert that there are no missing values

# A brief first conclusion

The predictor variables are all numeric

There are no missing values

The qualitative variables have already been treated ( no need to apply some technique like One hot Encoding, Ordinal Encoder ... )

There is almost no correlation between the predictive variables and the target variable (shares)

It looks like the regression is going to be complicated

The background is a dark, abstract composition. It features a grid of glowing blue and green binary digits (0s and 1s) scattered across the frame. Overlaid on this are several semi-transparent financial charts. On the left, there's a candlestick chart with red and green bars. In the center and right, there are line graphs with jagged, fluctuating lines in shades of blue and green. The overall aesthetic is high-tech and data-driven.

# Regression



# Regression

---

We'll start by using 2 different types of algorithms

- Linear Regression
- Random Forest Regressor

We will quickly observe the score we obtain with these algorithms before continuing or seeking to optimize the hyperparameters.

# Regression – Feature selection

---

- Even if we are not going to proceed to the hyperparameter tuning right away, we can already try to reduce the number of variables from which we will be working on.
- We use RFECV to do this. ( Recursive Feature Elimination Cross Validation )
- It returns us a list of parameters that we will keep and can even give us a ranking of our parameters

# Regression – Modeling

---

- As we don't obtain any good result in regression, we can try with a **classification objective**.
- We obtain a score of 0,11 in the best case with a Random Forest Regressor estimator.
- Now we will try to **determine if an article is popular** instead of the number of shares.



The background of the slide is a dark blue grid. Each cell in the grid contains a two-digit hexadecimal character (0-9, A-F) in a lighter blue color. The characters are arranged in a way that creates a sense of depth and digital data. The word 'Classification' is centered in the middle of the slide in a white, sans-serif font. A thin white vertical line is positioned to the left of the text.

# Classification



# Classification



AS WE CANNOT PRECISELY  
PREDICT THE NUMBER OF  
SHARES



OUR NEW GOAL IS TO PREDICT  
WHETHER THE ARTICLE IS  
POPULAR OR NOT.



Classification

But what does  
popular mean ?



# Classification

- That's why we create the variable 'popular' as a new target variable for a classification problem.
- The popularity threshold is defined by the median of the target variable.
- if the number of shares is greater than this value, we consider the article to be popular, otherwise, it is not.

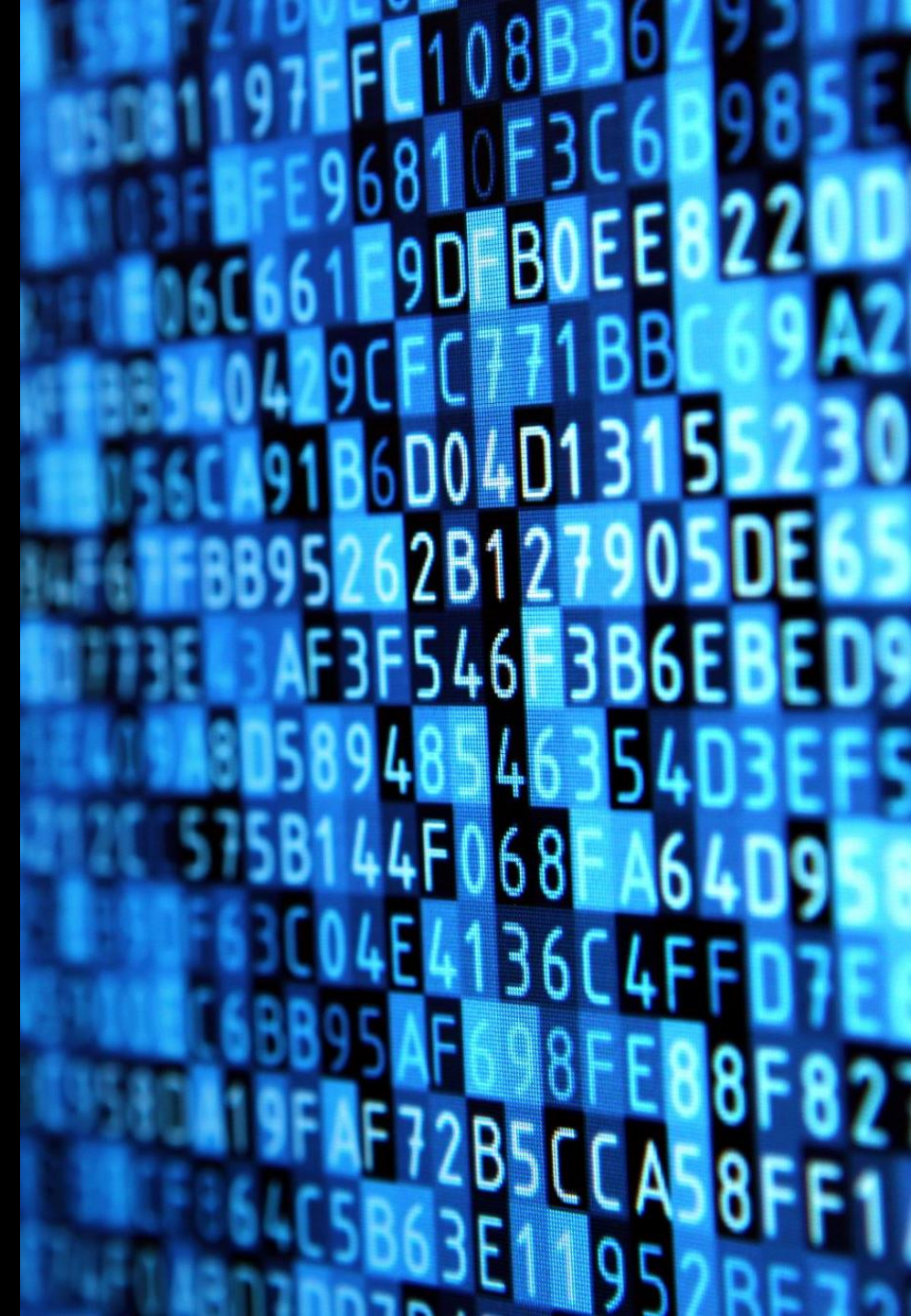


# Classification

We'll start by using 2 different types of algorithms

- SGDClassifier
- Random Forest Classifier

We will also quickly observe the score we obtain with these algorithms before continuing or seeking to optimize the hyperparameters.



# Classification – Feature selection

- 
- As with regression, we use RFECV to reduce the amount of parameter to keep.
  - However, we realize that the variables have almost no correlation with the target variable so we can hardly eliminate a lot of the variable.

# Classification – Modeling

- We obtain far better result than before.
- We obtain at least a score of 0,6 with SGDClassifier and 0,65 with a Random Forest Classifier.

That's why we decided to keep the Random Forest Classifier.



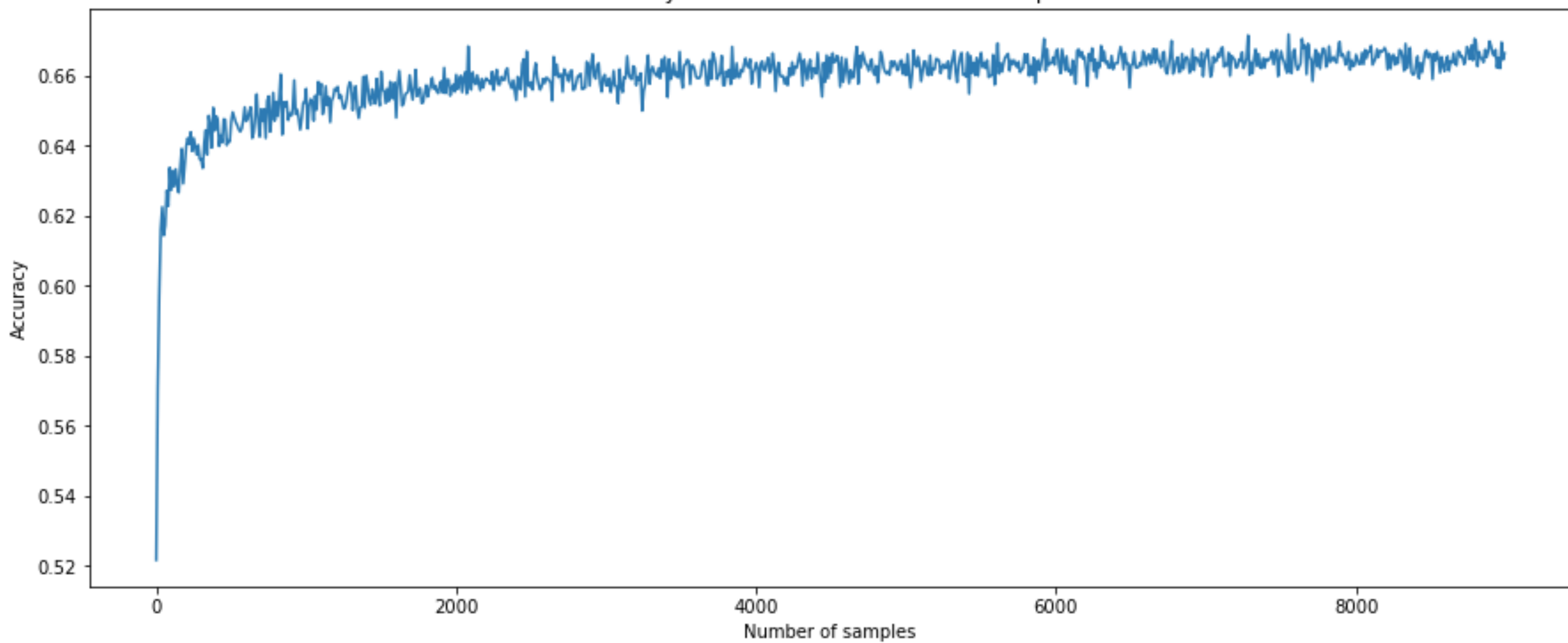
# Classification – Hyperparameter Tuning

For the tuning of the hyperparameters we used :

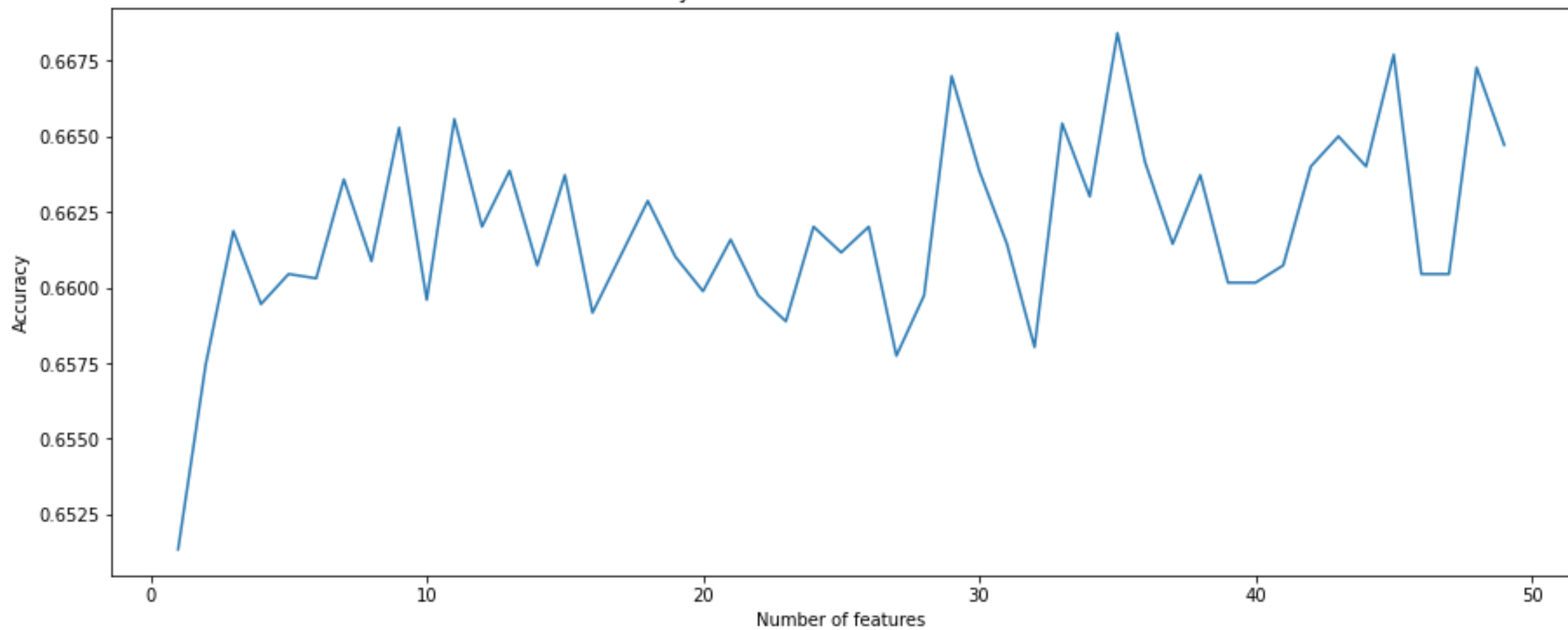
- Grid Search CV .
- We draw 3 learning curve.



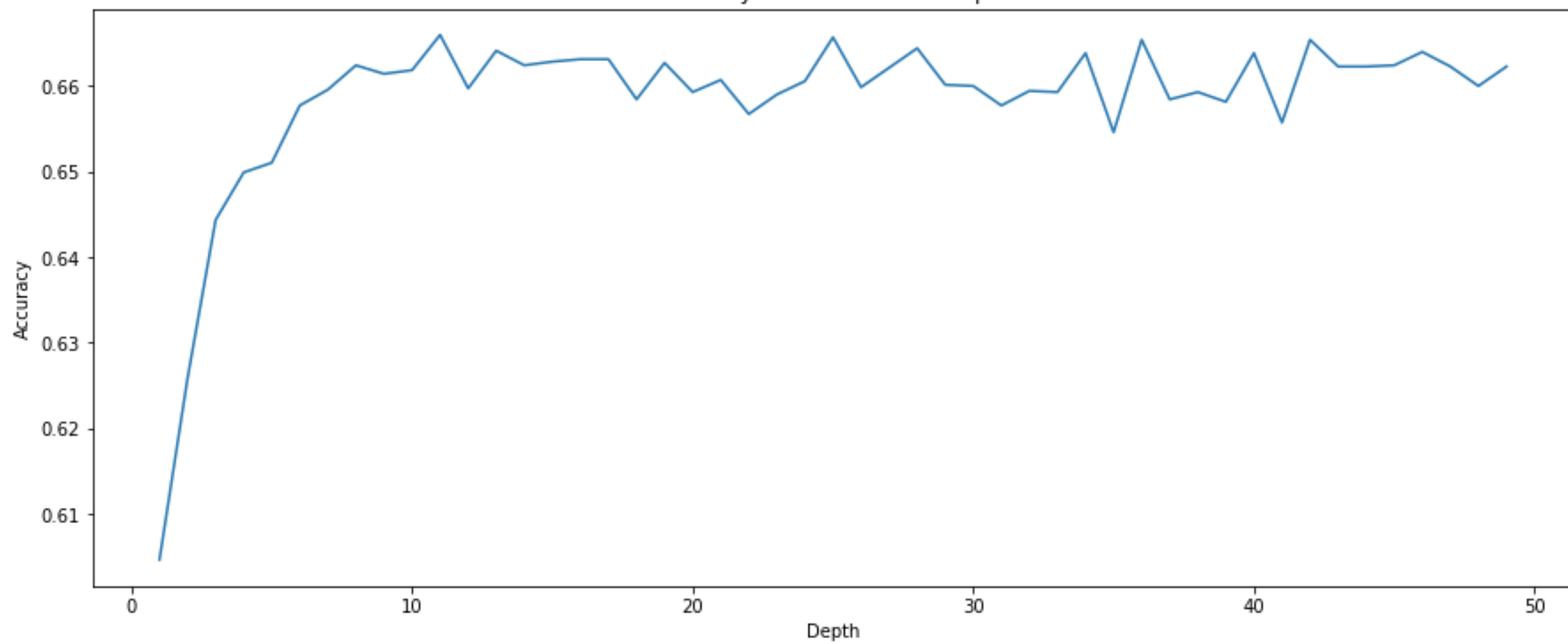
Accuracy in fonction of the number of samples



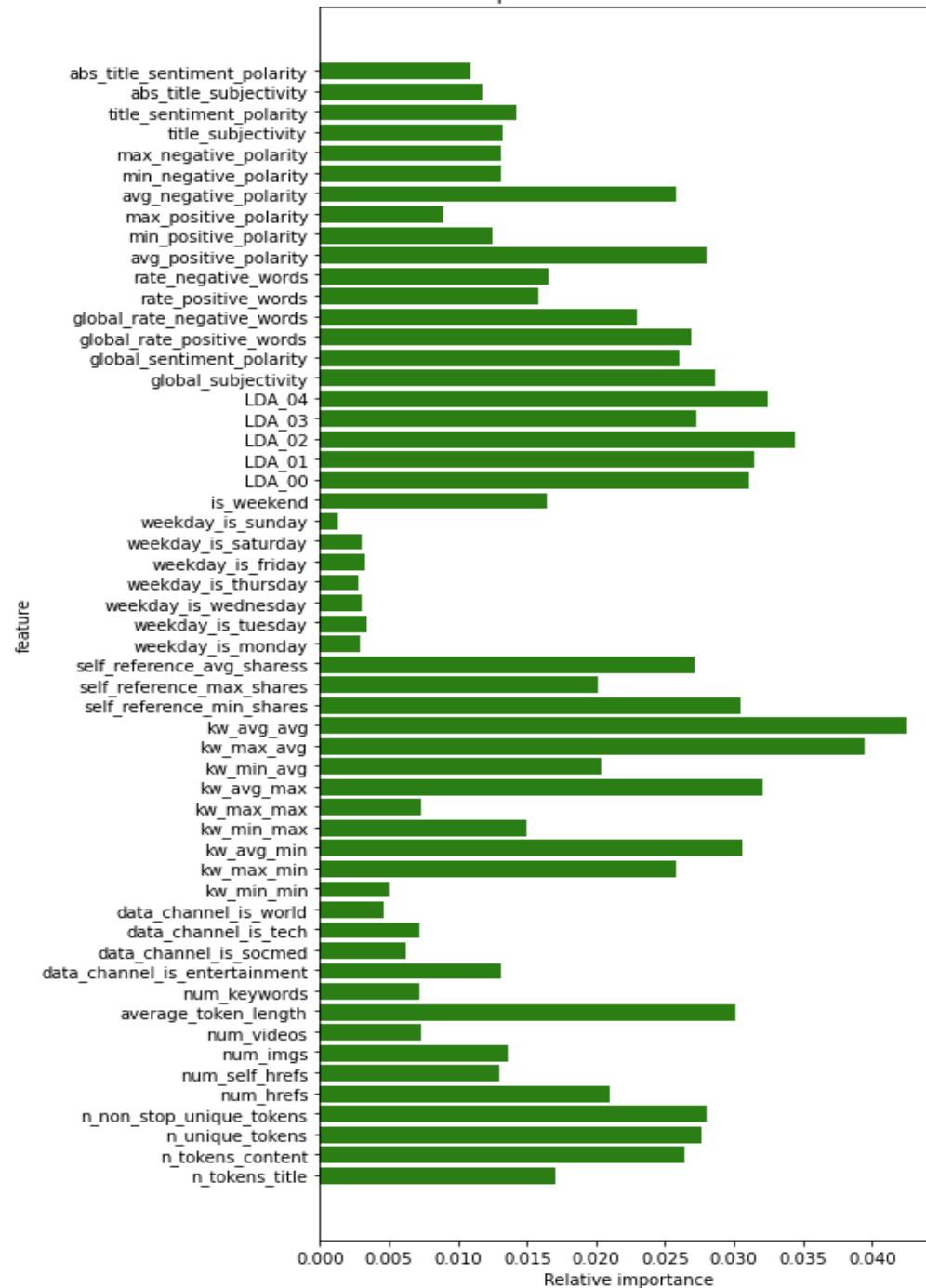
Accuracy in fonction of the number of features

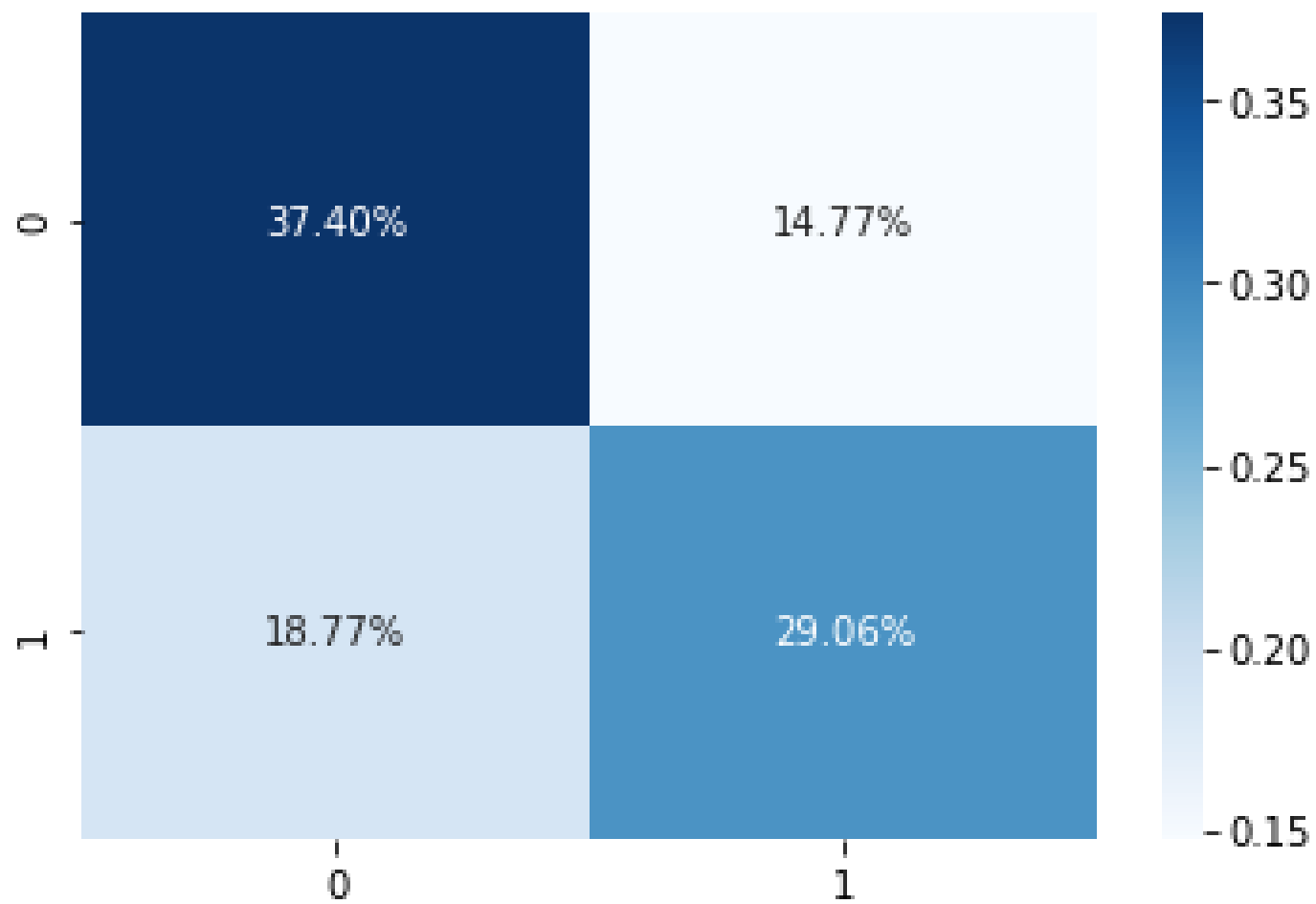


Accuracy in fonction of the depth



Feature importance in RandomForest Classifier





Confusion  
Matrix



A pile of 3D question marks of various shades of gray and white, scattered on a dark, textured surface. The lighting creates highlights and shadows, giving the question marks a three-dimensional appearance.

# Conclusion

- At the beginning the objective was to predict the number of shares but it appears that despite the number of data, the predictive variables defined do not allow us to solve the problem of regression.
- By defining the popularity of an article, we can redefine our problem to a classification problem that we can answer with our data.