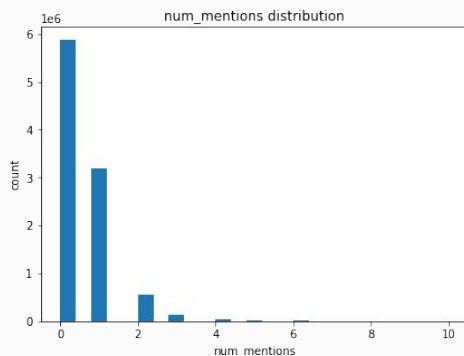# Data Mining project

A.Y. 22/23

Group 12
Leonardo Vona

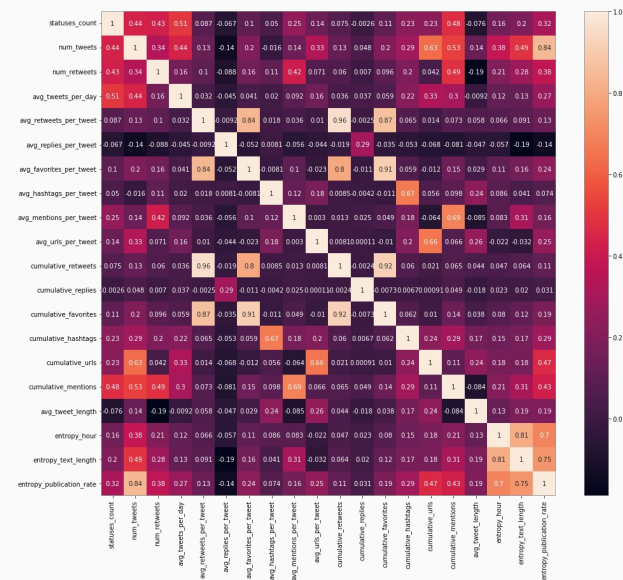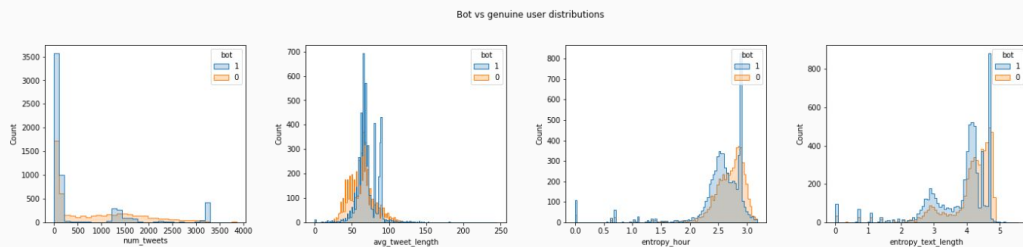# Data understanding & preparation

- Most of the features have a very high skewed distribution.



- Two tweets were considered equal if they share the same value for *user_id*, *created_at*, and *text*

- The features were checked for correctness by also using external sources

- For most cases, the missing values have been substituted after data cleaning with the median

- The tweets data frame has been separated into tweets and retweets

# Data understanding & preparation: new features

The new features have been extracted by joining the users and tweets data frames. For each user, indicators were extracted about its behavior in terms of activity and entropy.



Bot vs genuine user distributions
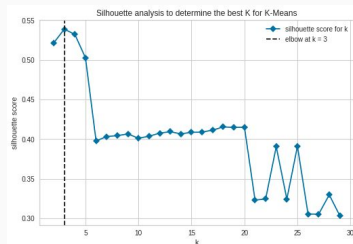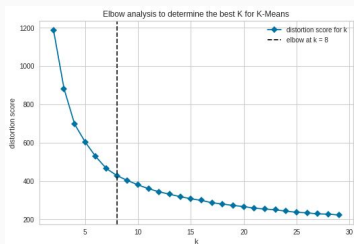
# Clustering

**Clustering techniques**

- KMeans
- Density-based
- Hierarchical
- XMeans (pyclustering)
- BSAS (pyclustering)

**Preprocessing**

- Extract numerical features from the dataset describing the user behaviour
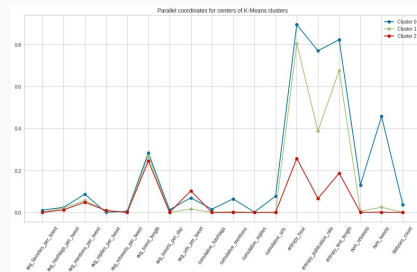- Remove highly correlated features
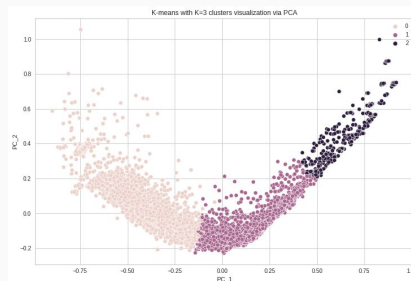- Apply normalization

# Clustering: KMeans

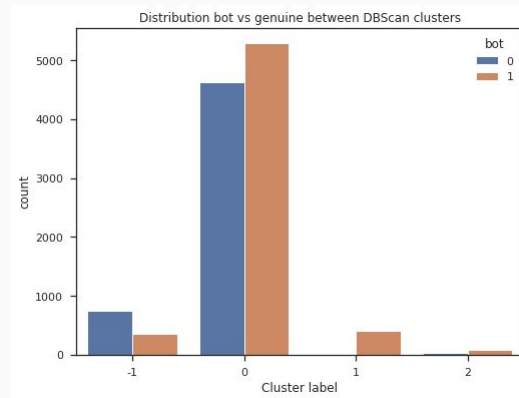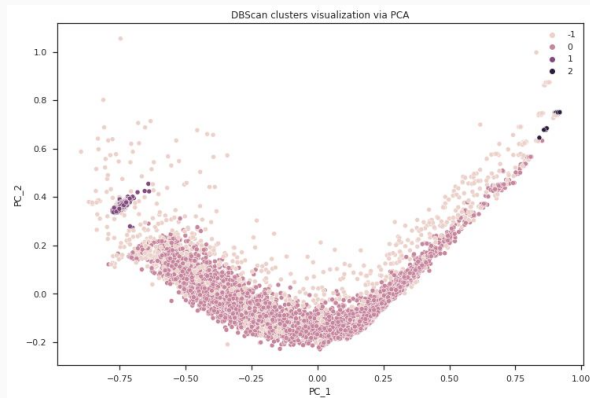The best K has been determined using the elbow and silhouette methods.



By using SSE, Davies-Bouldin, Silhouette and Calinski-Harabasz scores, the K value of 3 has been chosen from the candidates.

- Cluster 0 groups the most active users
- Cluster 1 contains medium active users
- Cluster 2 contains users with a scarce activity
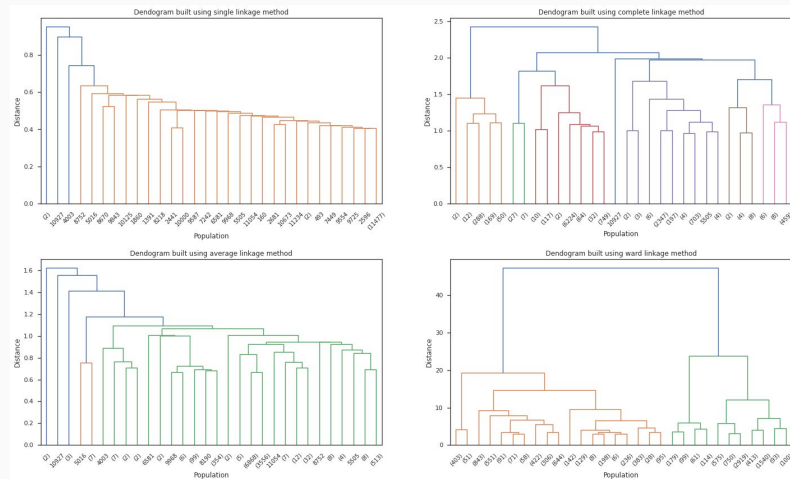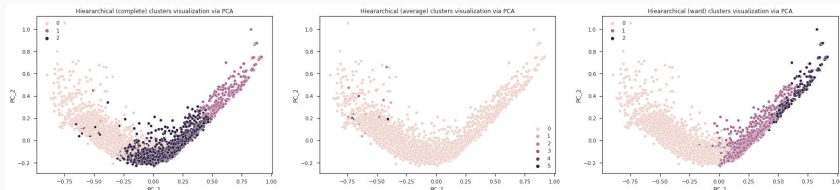
# Clustering: Density-based

A grid search is performed to establish the best combination of *minPts* and *eps* parameters.





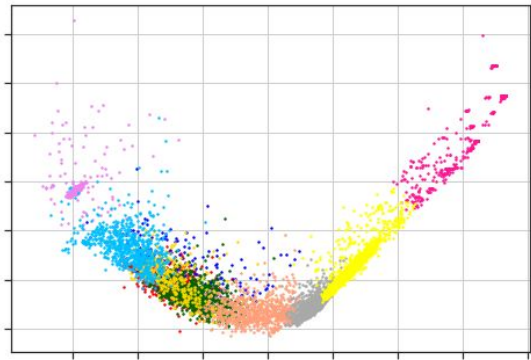Density based clustering is not particularly suited for our dataset.

# Clustering: Hierarchical

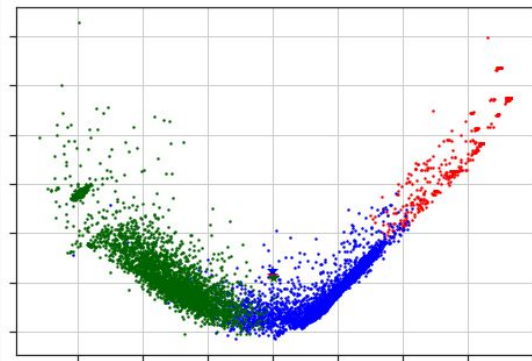The Complete and Ward linkage methods seem the best options.

# Clustering: pyclustering



X-Means clustering tends to split the data in the maximum number of clusters, because it gives the best BIC score.

BSAS clustering makes a good separation of the data, and is particularly suited when the data is given as a stream.

# Classification

**Preprocessing**

- Drop *id* and *name*, exclude *bot* attribute
- Group small sized languages together
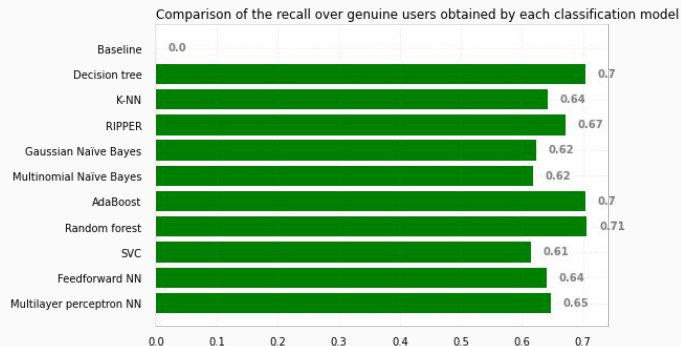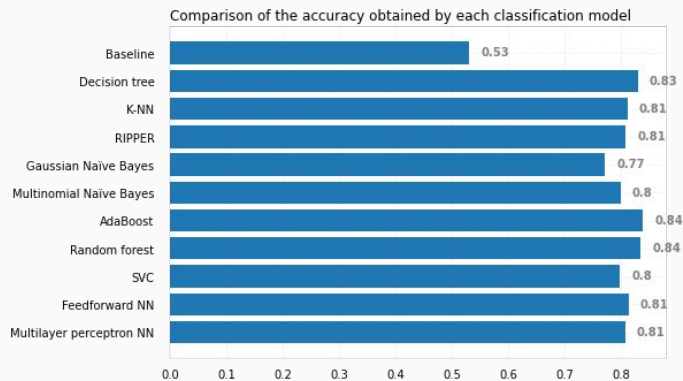- Extract new features from *created_at*

**Models**
Decision Tree, K-NN, Rule-based, Naive Bayes, AdaBoost, Random Forest, SVM, Neural Network

**Approach**

- Apply specific preprocessing for the model (e.g. normalization, reduce dimensionality, convert features)
- Split in train and test set
- Grid (or randomized) search with CV to select hyperparameters
- Train and evaluate

# Classification: comparison

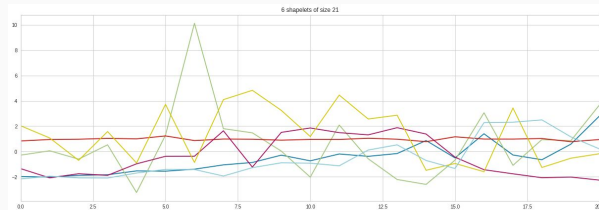The best classifiers are the Ensemble Methods, reaching an accuracy of 84%.

Comparison of the accuracy obtained by each classification model

| Model | Accuracy |
|---|---|
| Baseline | 0.53 |
| Decision tree | 0.83 |
| K-NN | 0.81 |
| RIPPER | 0.81 |
| Gaussian Naïve Bayes | 0.77 |
| Multinomial Naïve Bayes | 0.8 |
| AdaBoost | 0.84 |
| Random forest | 0.84 |
| SVC | 0.8 |
| Feedforward NN | 0.81 |
| Multilayer perceptron NN | 0.81 |

Comparison of the recall over genuine users obtained by each classification model

| Model | Recall |
|---|---|
| Baseline | 0.0 |
| Decision tree | 0.7 |
| K-NN | 0.64 |
| RIPPER | 0.67 |
| Gaussian Naïve Bayes | 0.62 |
| Multinomial Naïve Bayes | 0.62 |
| AdaBoost | 0.7 |
| Random forest | 0.71 |
| SVC | 0.61 |
| Feedforward NN | 0.64 |
| Multilayer perceptron NN | 0.65 |

All the models are unable to discriminate well the genuine users, with a general high rate of false positives.
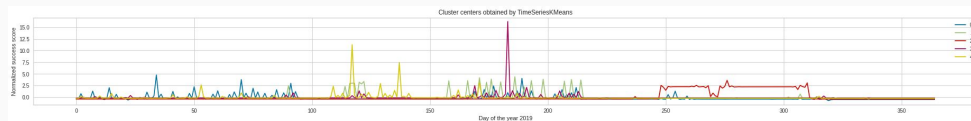
# Time series analysis

## Preprocessing

- Compute Success Scores of 2019 for each user
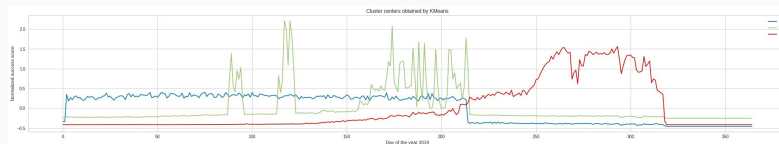- Fill missing values
- Apply amplitude scaling

## Shapelet extraction



## Clustering

- Partitional



- Feature-based

Thanks!