

Modern Data Analysis

Aleksey Ryabykin, Mikhail Osokin

20/12/22



Datasets

We have chosen a [dataset](#) consisting of Medium articles data. Our objective is to examine the topics, volume, content, and other factors that can contribute to the popularity of the articles on this platform.

Claps could be viewed as a platform business metric. Because anyone can clap as many times as they want, unlike with likes on other platforms, it's also intriguing to examine the outliers in this feature.

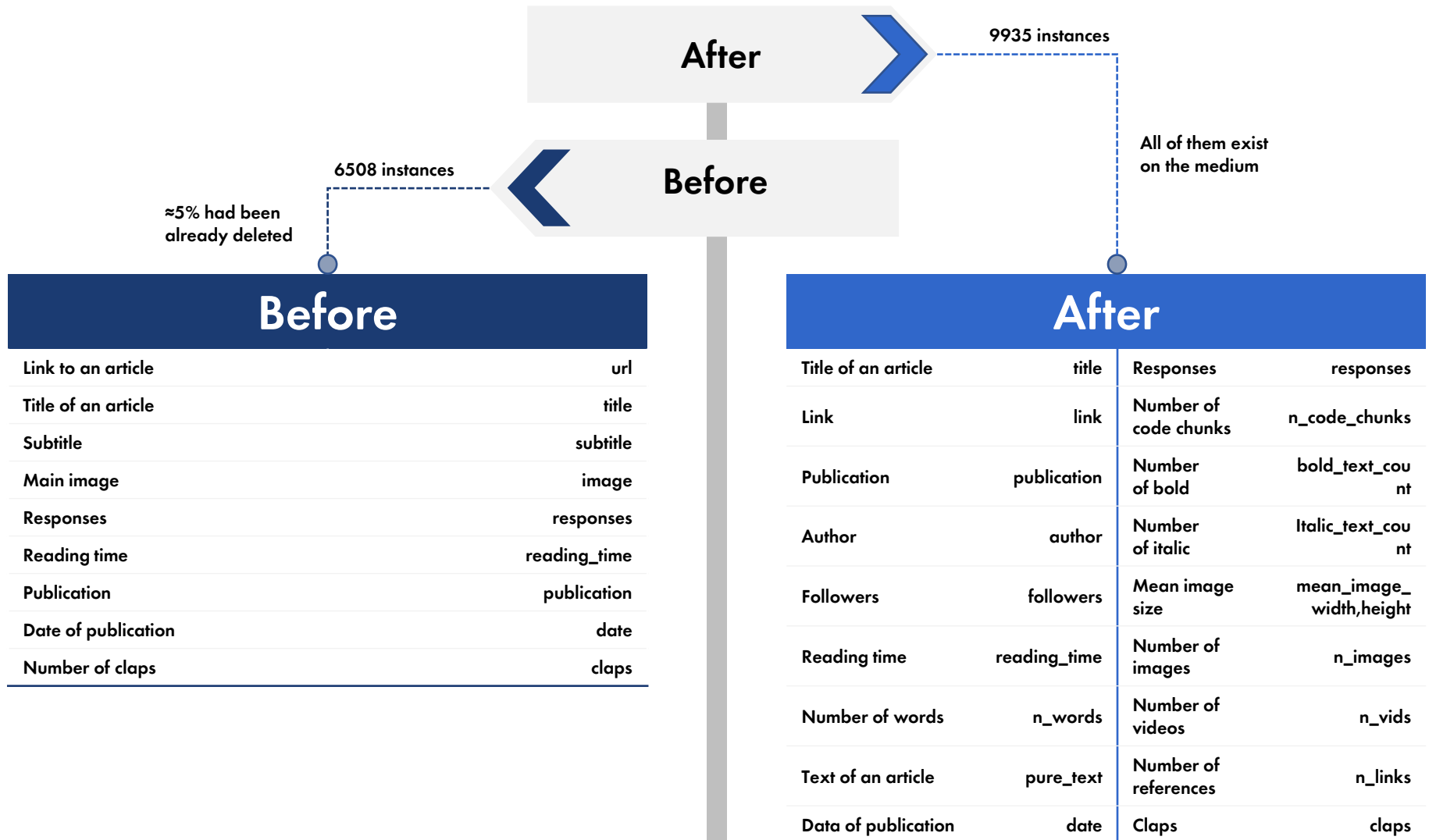
You also can read the variables description [here](#).



Proposed approaches

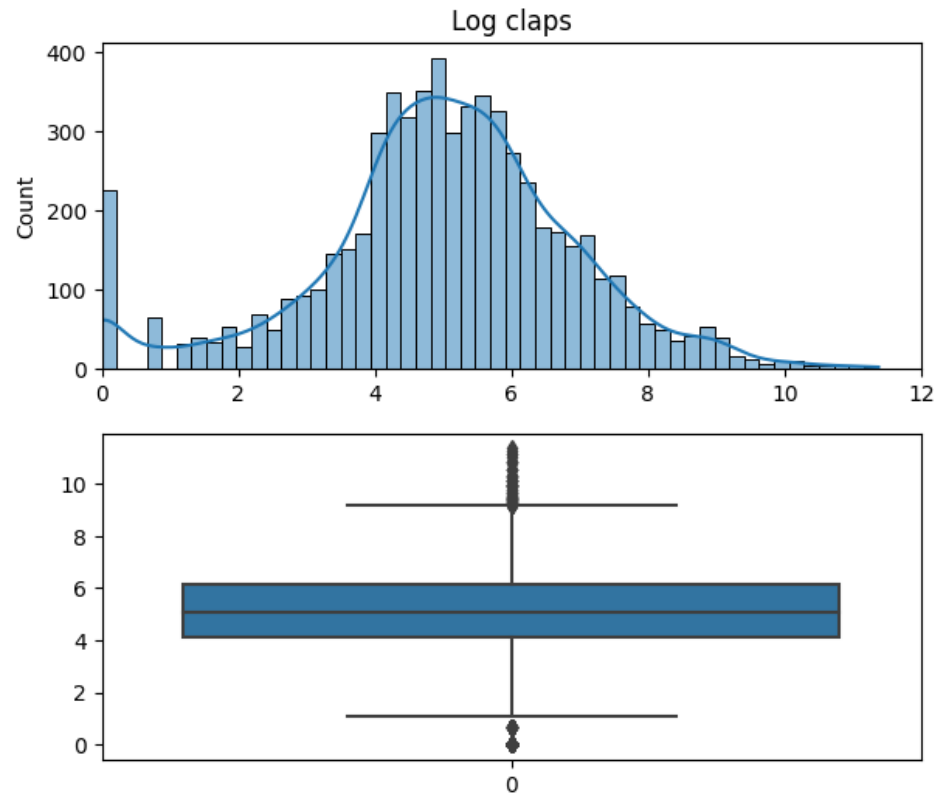
- Data scrapping (due to the date limits, we can scrap modern articles);
- Exploratory data analysis;
- Hypotheses testing;
- Topic modeling;
- Interpret the relation between topics and claps on the platform (if it exists);
- Build a regression with claps counts as a target variable;
- Build an ensemble model of regressors;
- Explain the results.

Data scrapping

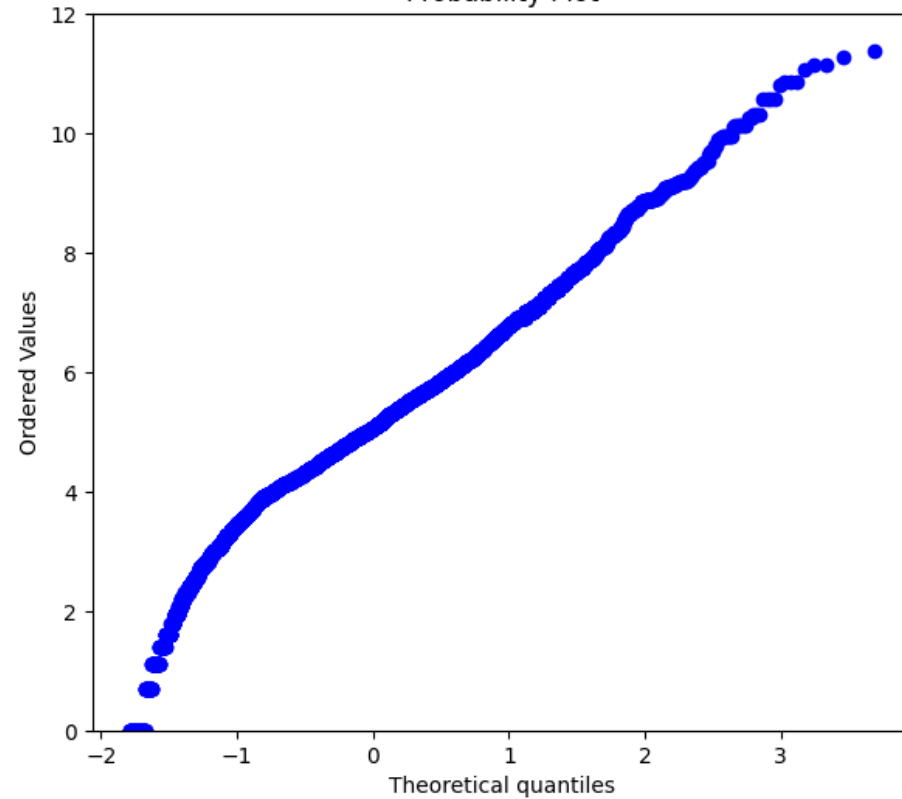


EDA inferences

Log claps variable



Probability Plot

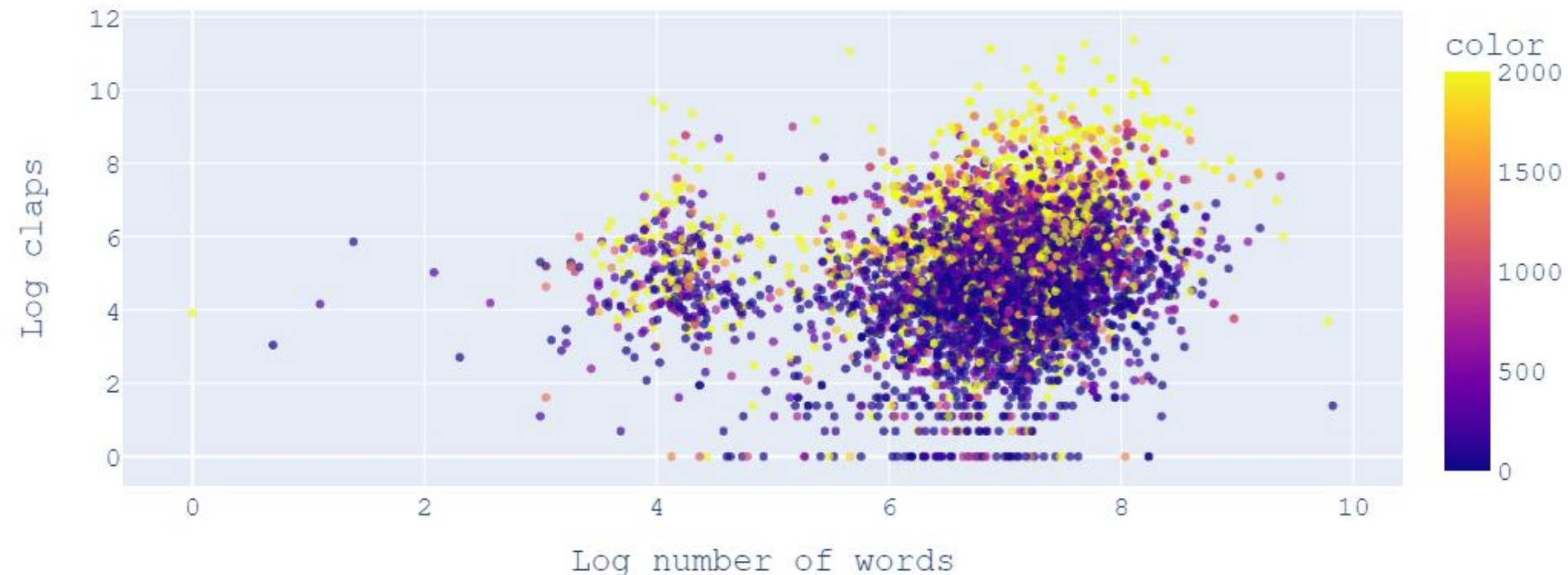


Close to normal distribution

EDA inferences

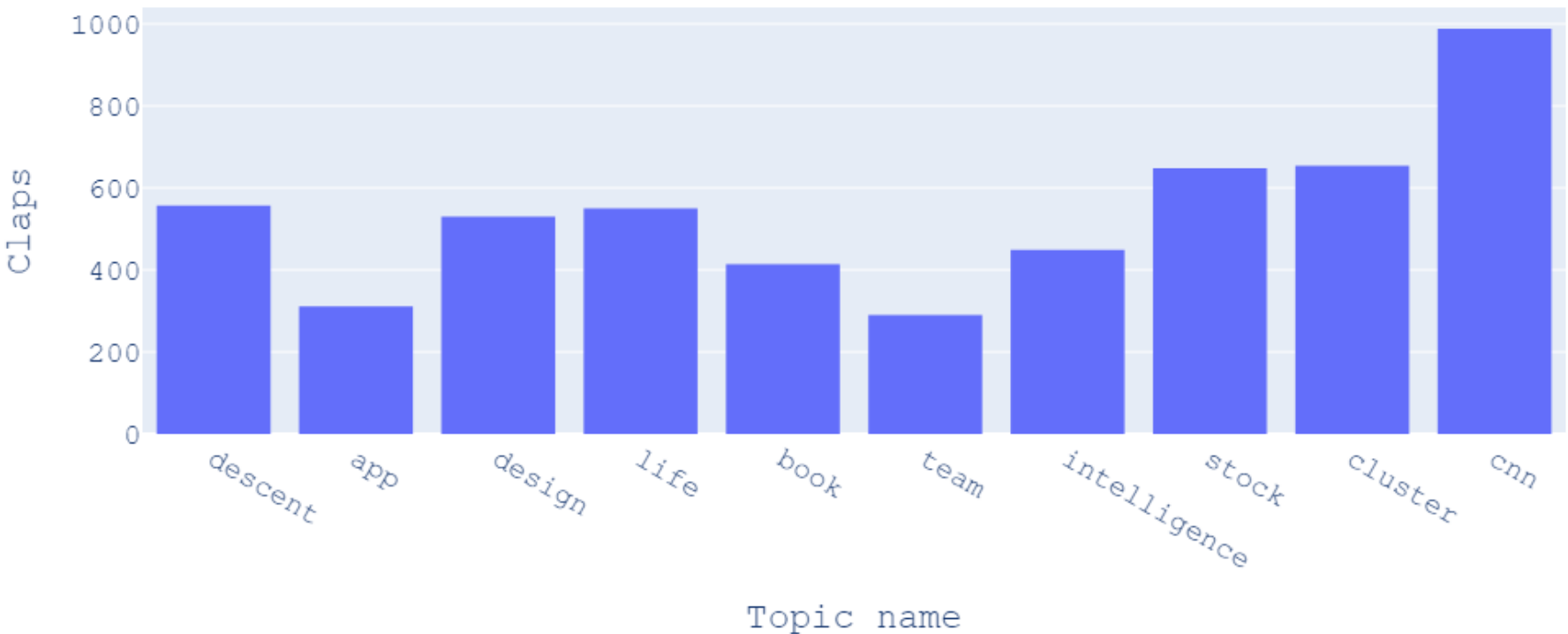
As we see, data splits into two clearly separated groups. Moreover, this plot shows that authors with more followers get more claps. However, there is no visible dependency of number of followers and length of articles.

Followers (as color), number of words and log claps



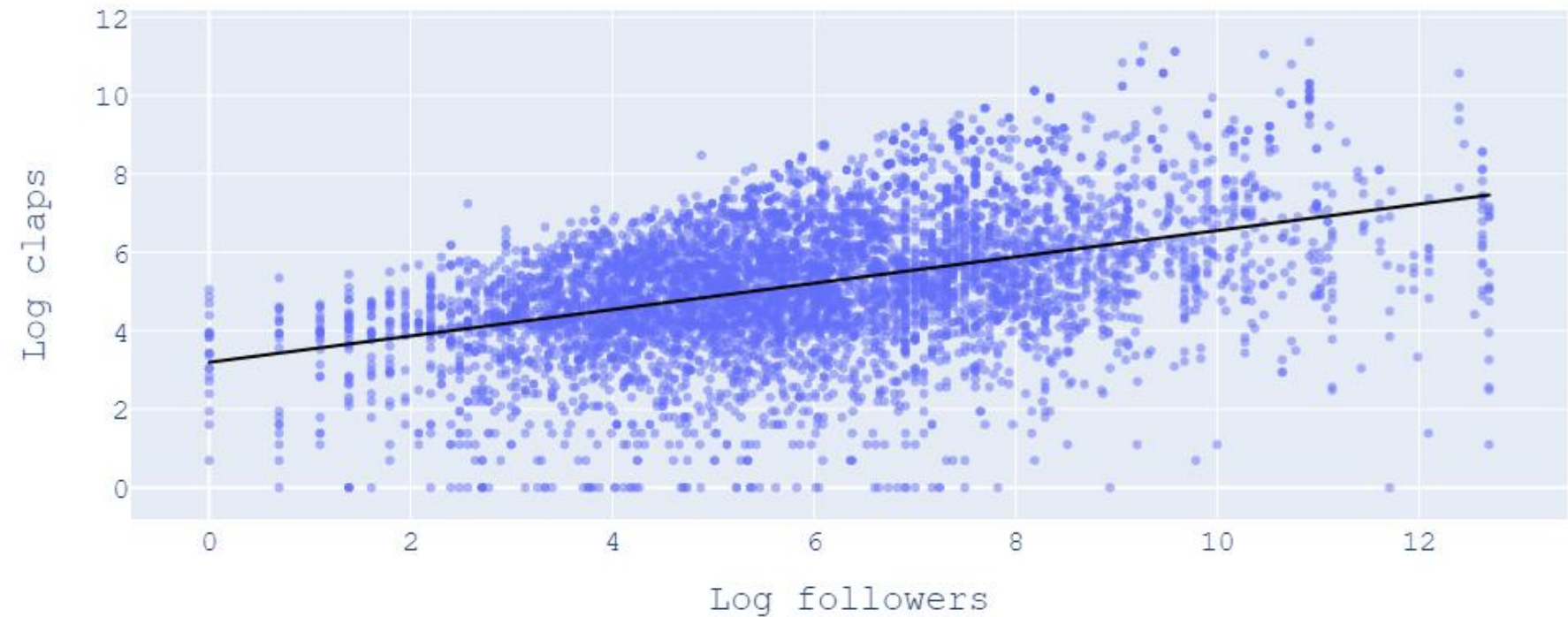
EDA inferences

The most popular topics



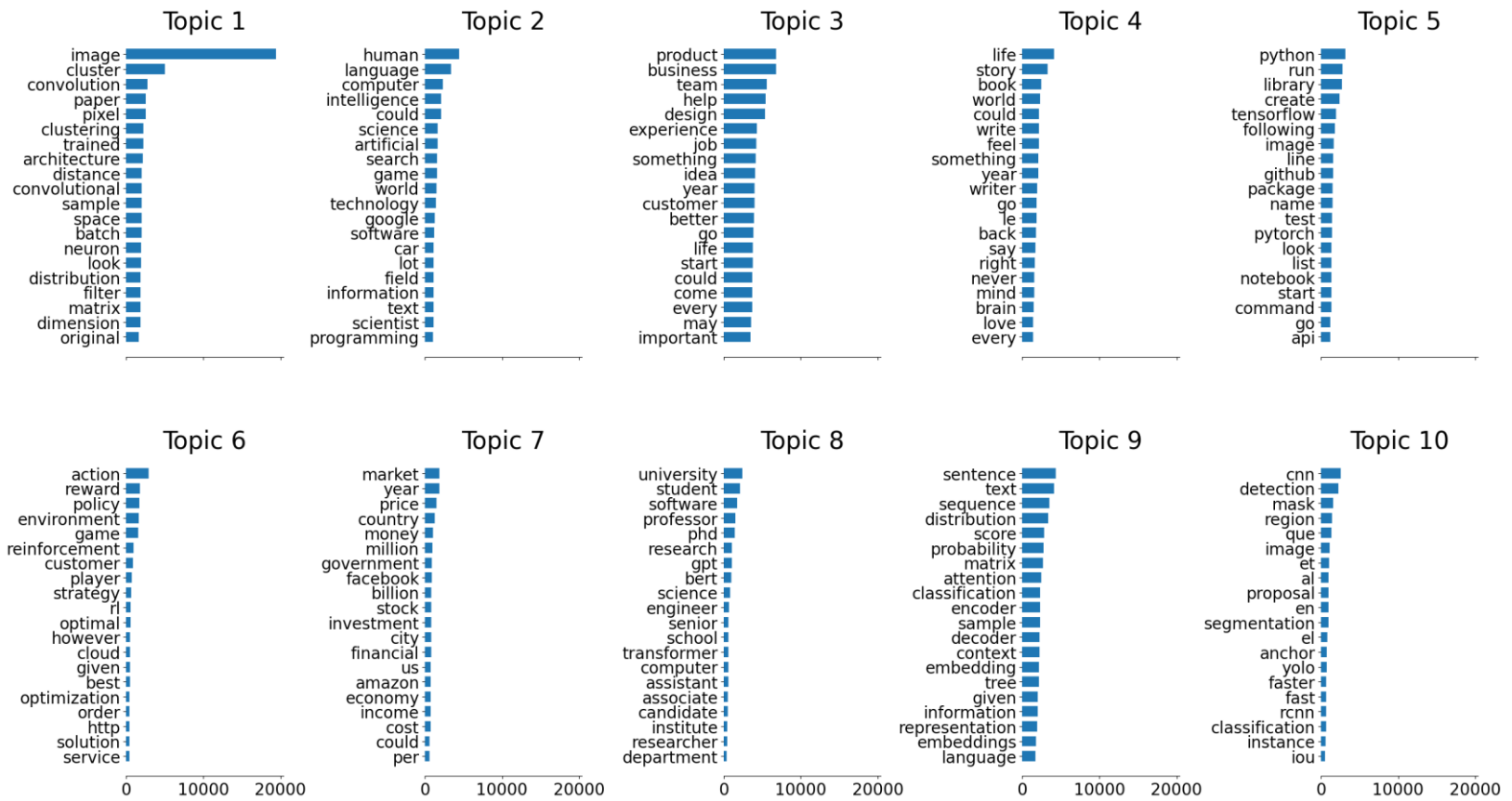
EDA inferences

Dependence between log followers and log claps

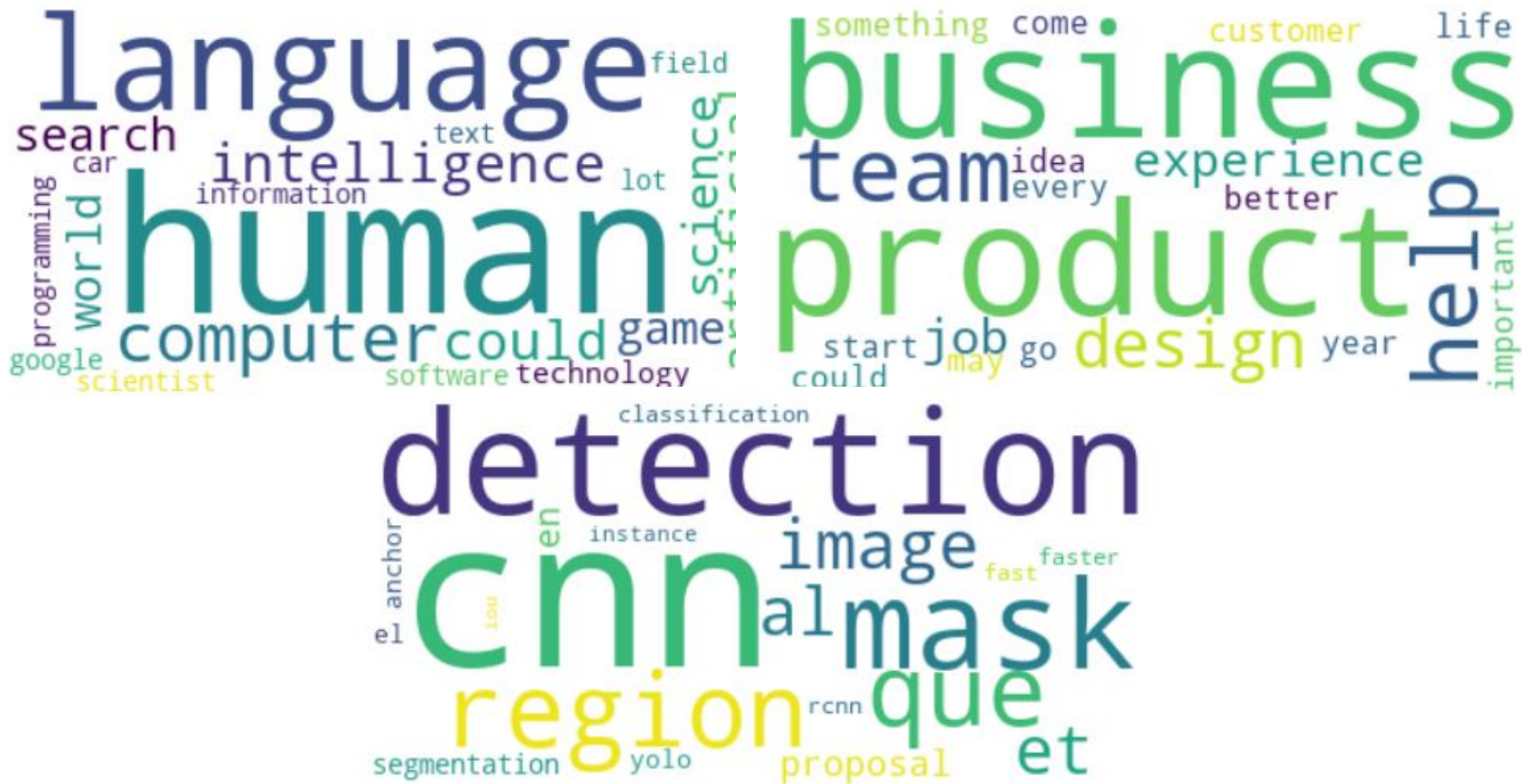


Topic modeling LDA

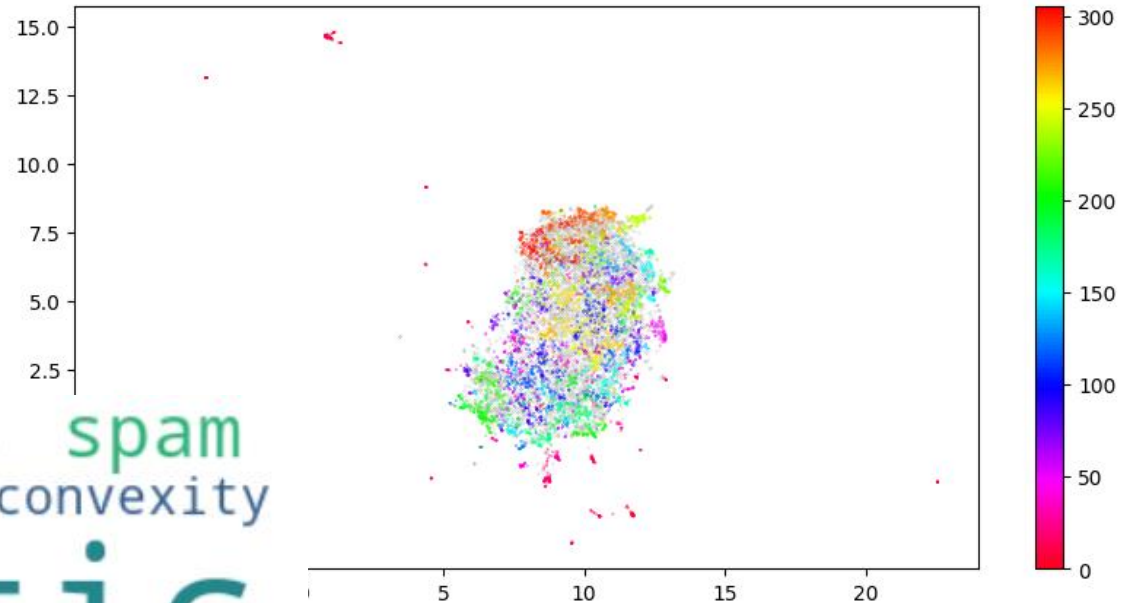
Topics in LDA



Topic modeling LDA



Topic modeling BertTopic (Hand-made)



meta lyapunov spam
convexity
logistic
descent
regression W convergence fraud

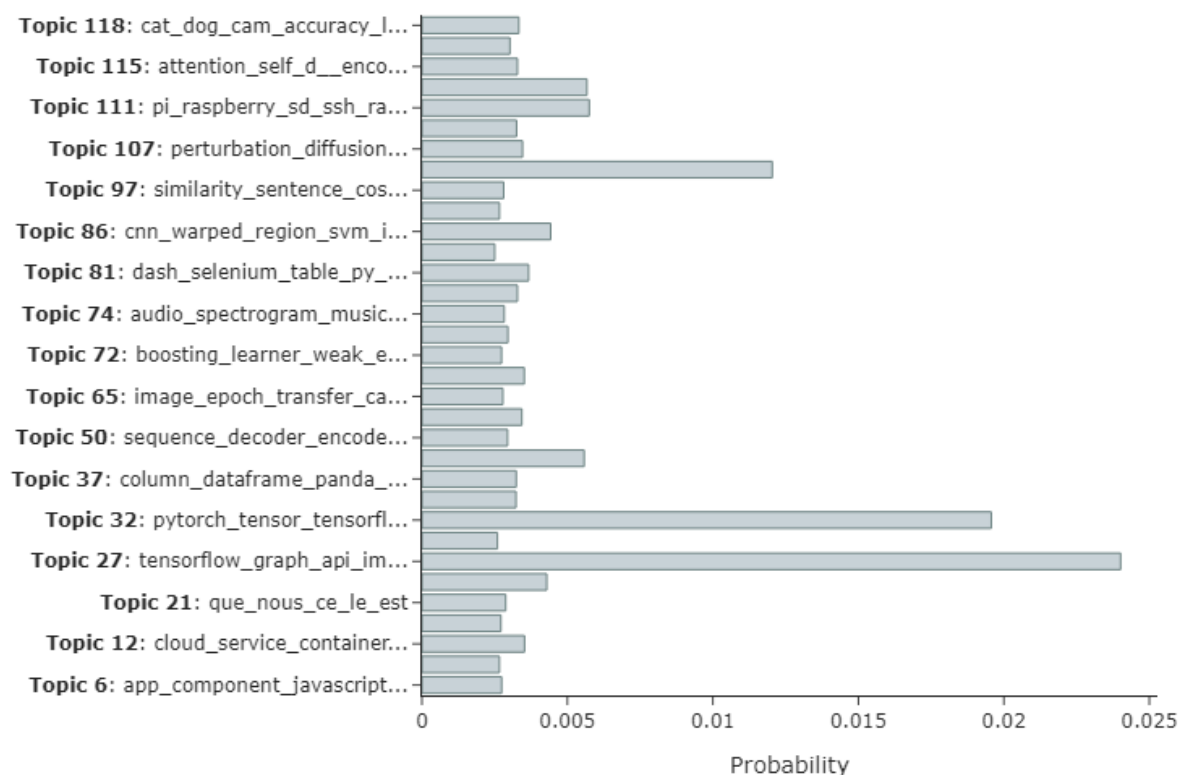
Topic modeling BertTopic

Topic Word Scores



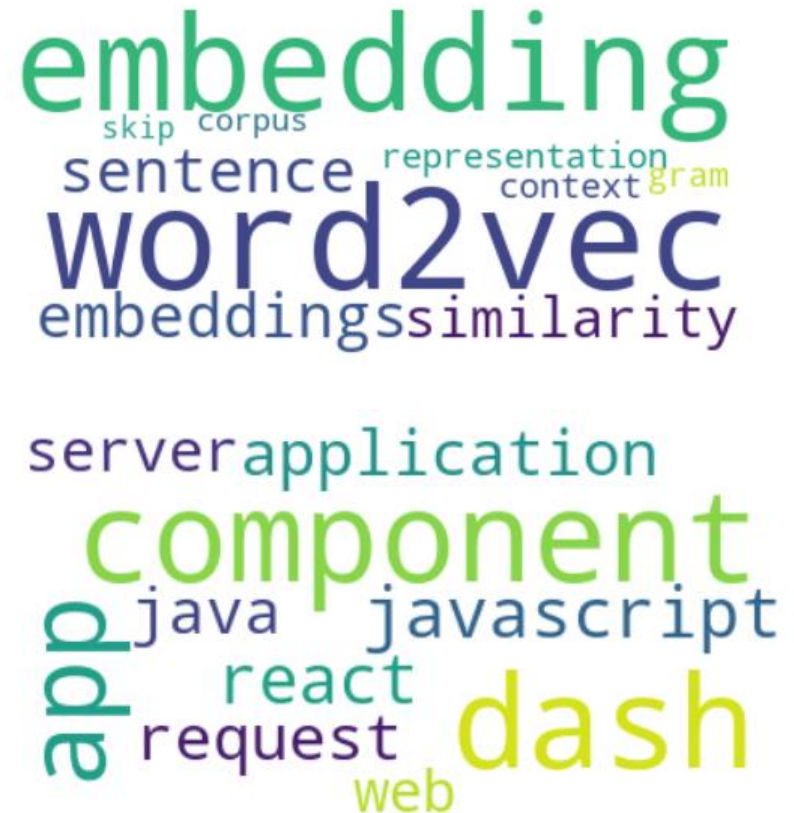
Topic modeling BertTopic

Topic Probability Distribution



Inferences of topic modeling

All clusters are easy to interpret and useful working with target. For instance, if we do not take the most popular ones, then the variance and expectation are preserved for the same clusters. It indicates a fairly close relationship.



Regression results

loss	catboost	Xgboost	Lightgbm
mae	307.4	297.3	353
mse	371.67	307.7	408.6
huber	453.2	450.1	476.4

Optuna-tuned ensemble of regressors

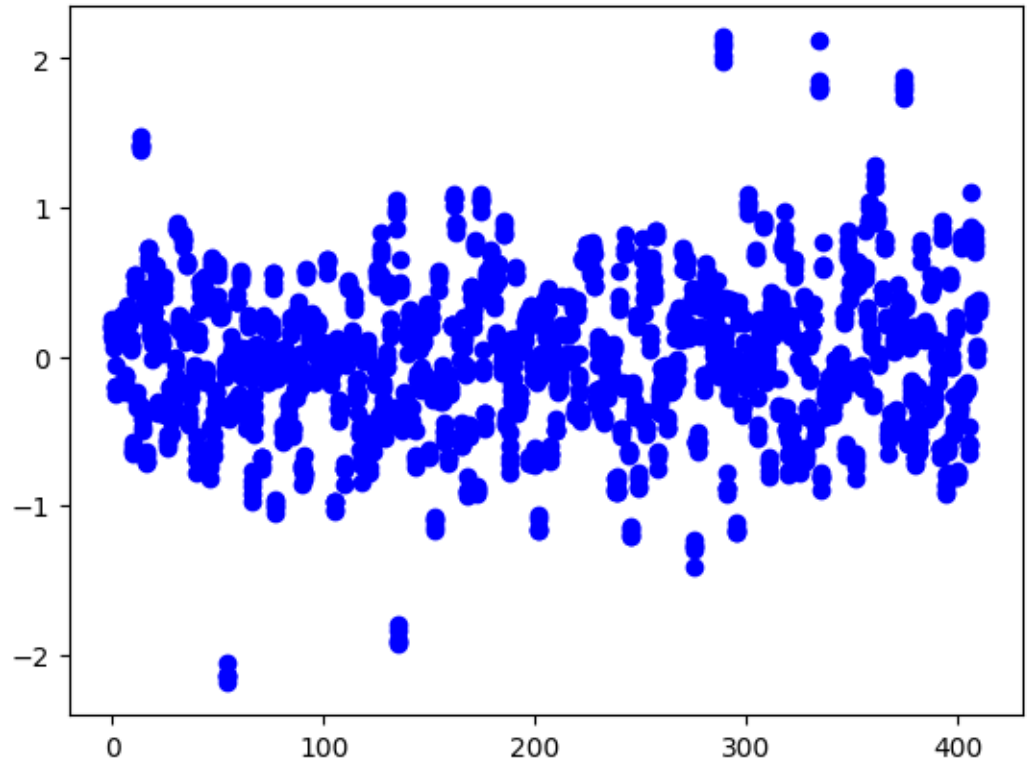
optuna	catboost	Xgboost	Lightgbm
282.2	307.4	297.3	353



OPTUNA

Error analysis

Regression residuals
proved their
homoscedasticity
following the Levene
criteria



	W	Pval	equal_var
Levene	0.008147	0.999869	True

Overview

Time improvement by
asynchronization

$\approx 5x$

Scrapped articles

4120

Time improvement

$\approx 6.5x$

Average quality improvement

$\approx 32\%$

Number of ensembled models

3

Hypotheses tests

10