

Telegram data-analysis

Student: Stoliaruk Yuliia
Teacher: Andrew Kurochkin

Plan of the presentation

1. Introduction
2. How to get data
3. Exploratory Data Analysis
4. Final results
5. Further work
6. Git Repository

Introduction

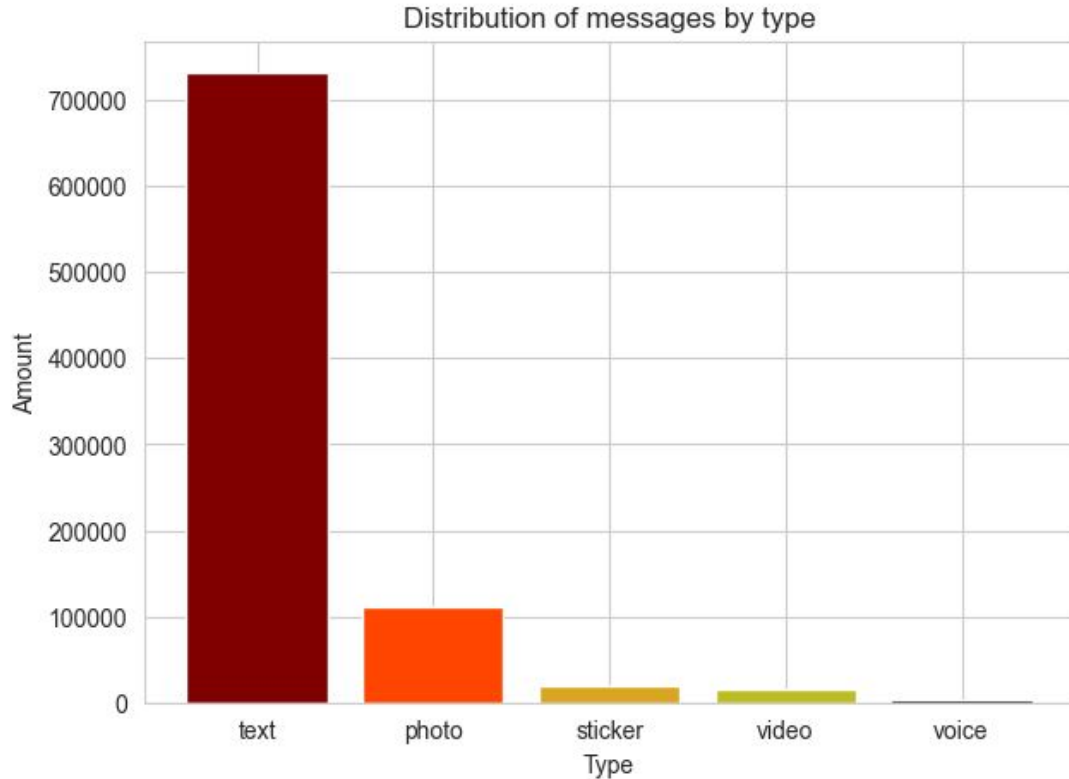
The main purpose of this project is to analyse activity of users and find hidden patterns in messenger Telegram.

How to get data

To get your data from Telegram:

1. To download data from Telegram use this repository [SanGreel/telegram-data-collection \(github.com\)](https://github.com/SanGreel/telegram-data-collection)
2. Then merge collected data using this repository [SanGreel/telegram-dialogs-analysis-v2 \(github.com\)](https://github.com/SanGreel/telegram-dialogs-analysis-v2)

Data overview (Messages data)

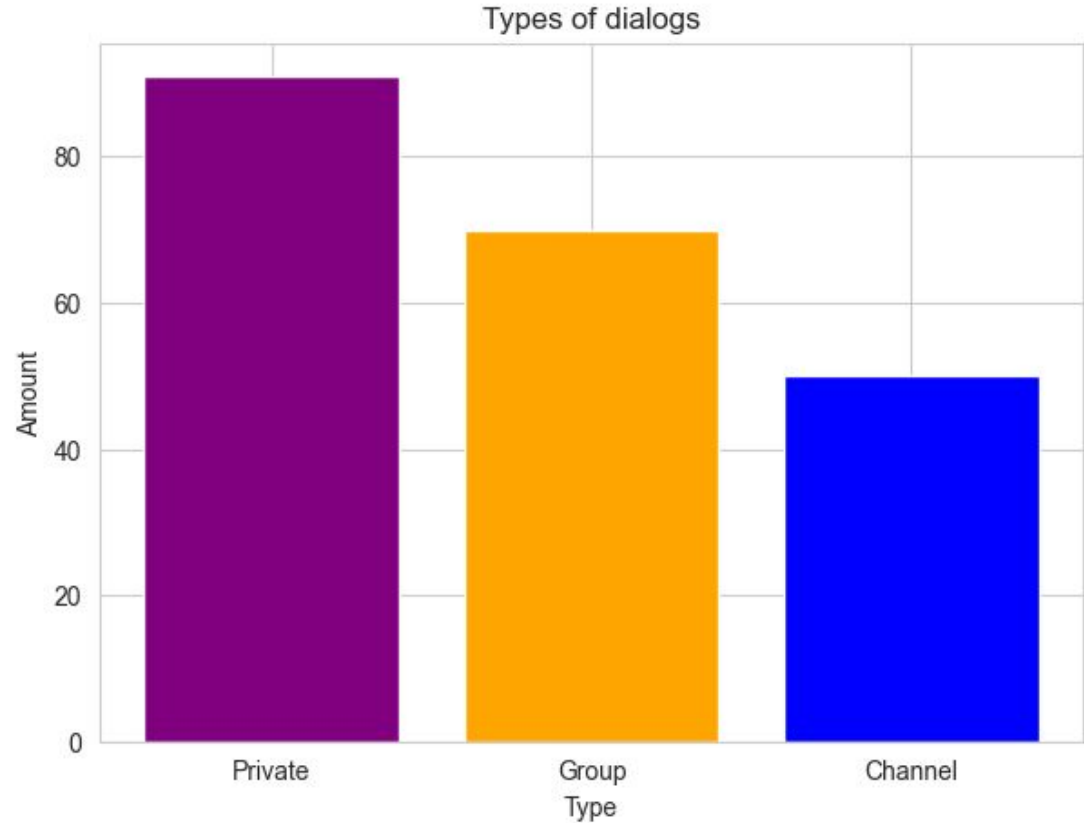


Size: 337.849 MB

Messages number: 880 940

Data overview (Dialogs data)

Size: 1.85 MB

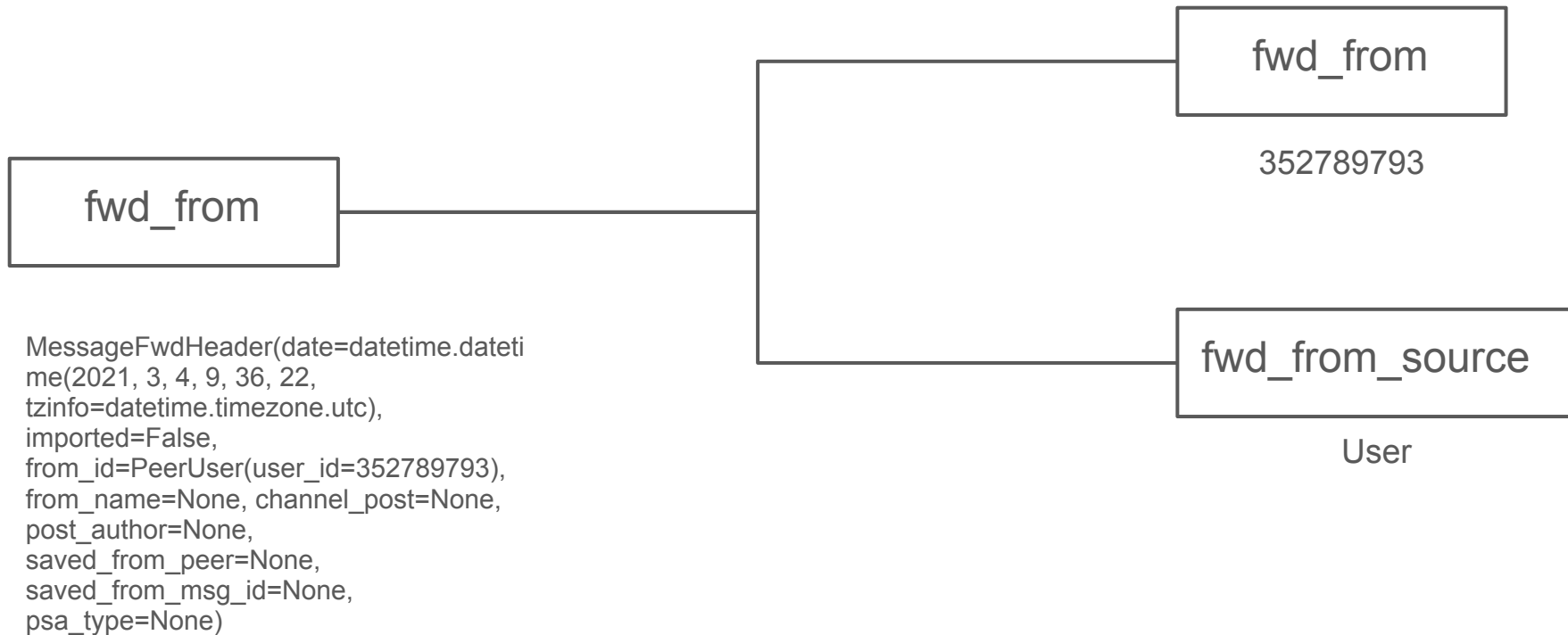


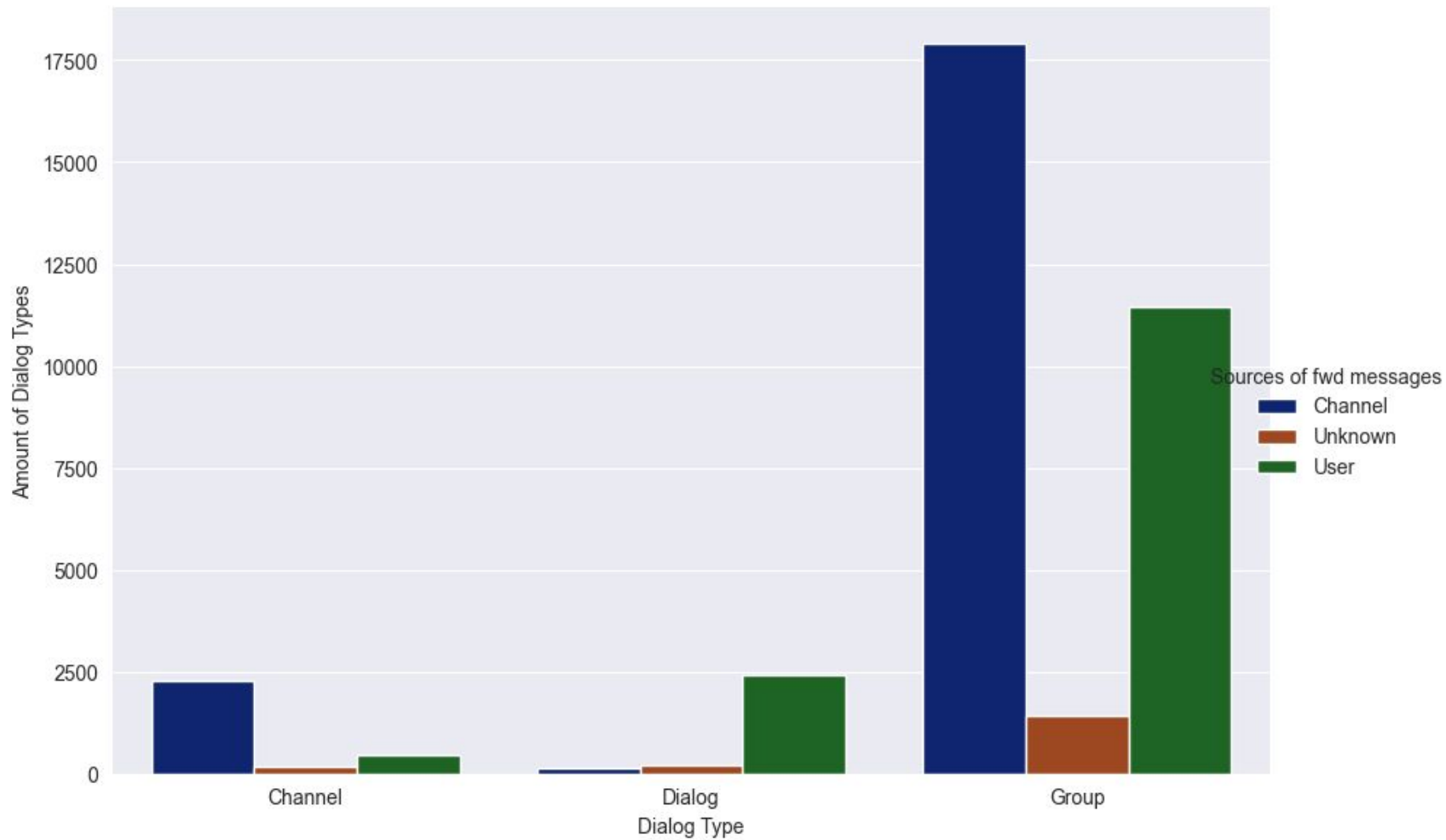
Exploratory Data Analysis (Data cleaning)

The first step was to clear the data.

1. Cleaning 'fwd_from' column
2. Cleaning 'from_id' column
3. Cleaning 'message' column

Cleaning 'fwd_from' column





Cleaning 'from_id' column

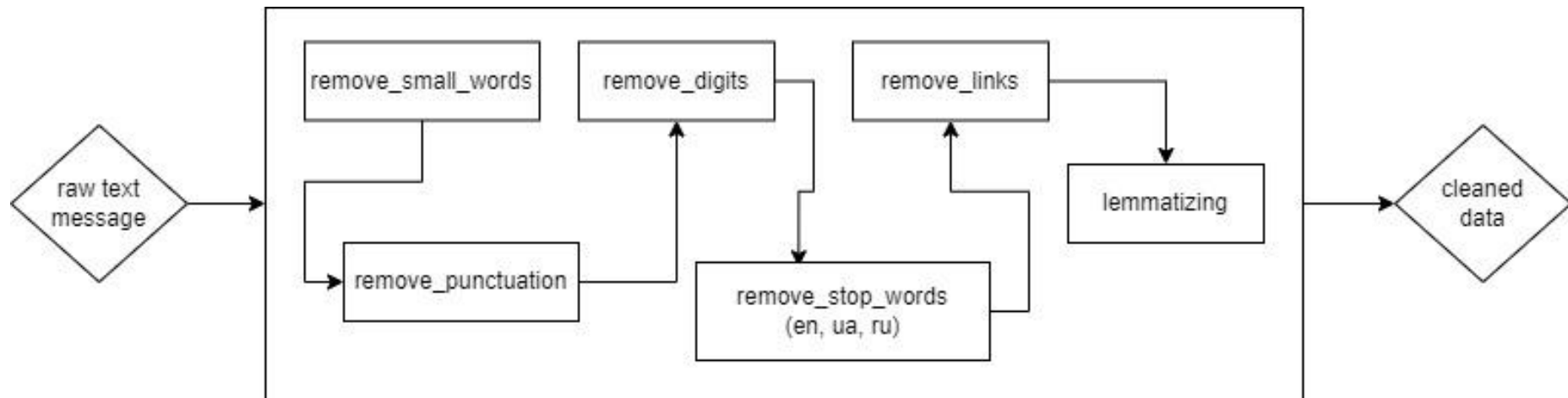
I cleaned this column for easier analysis.

`'PeerUser(user_id=475253228)'`



`'475253228'`

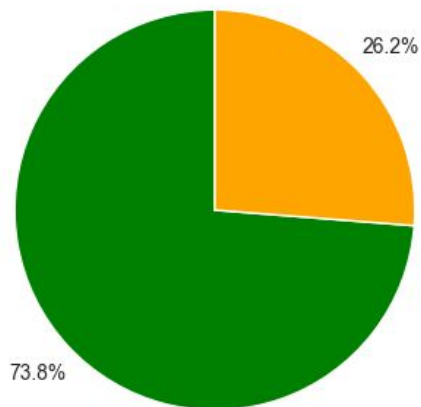
Cleaning 'message' column



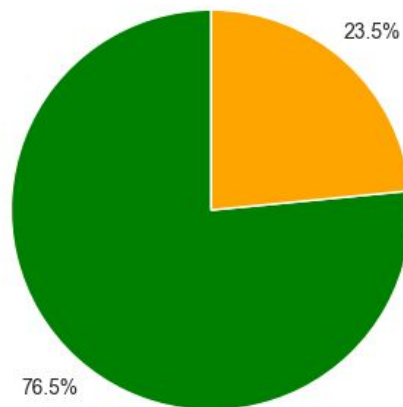
General analysis of my messages

Language distribution

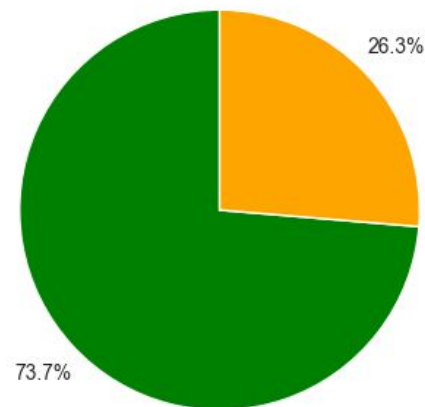
All Messages



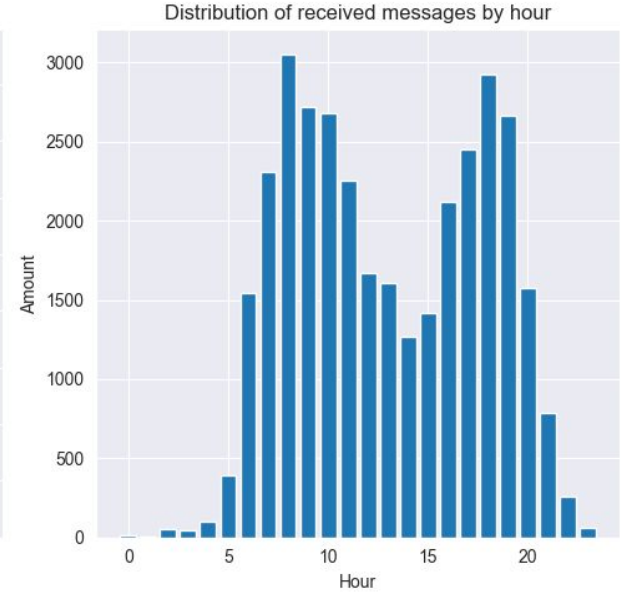
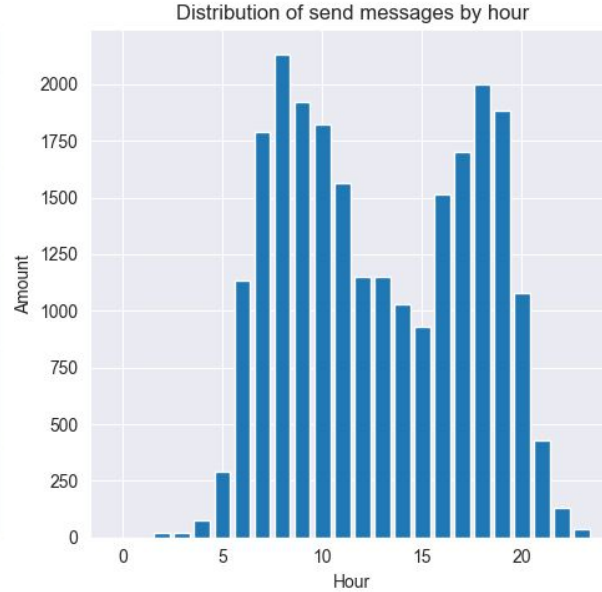
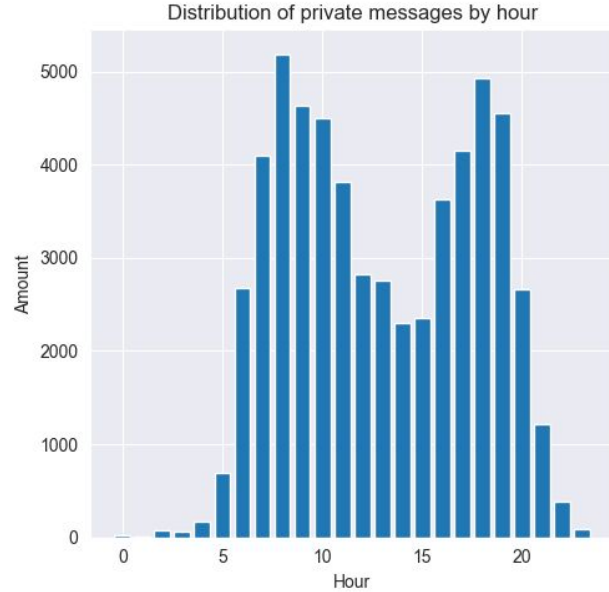
Sent Messages



Received Messages

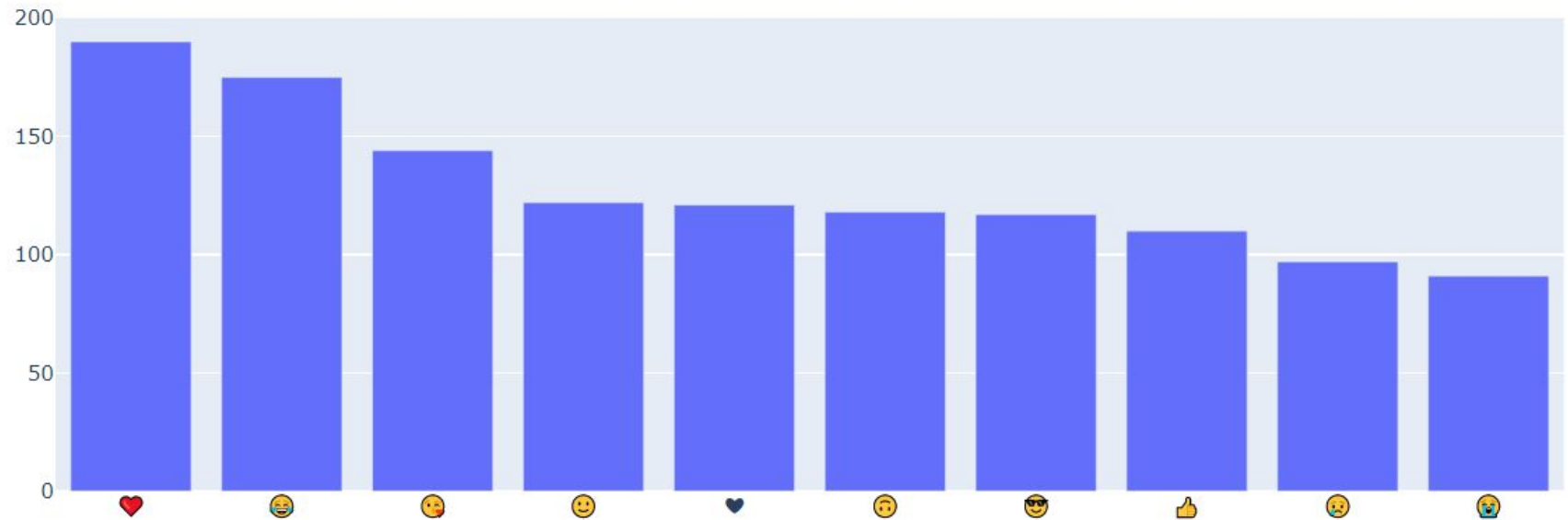


General analysis of my messages



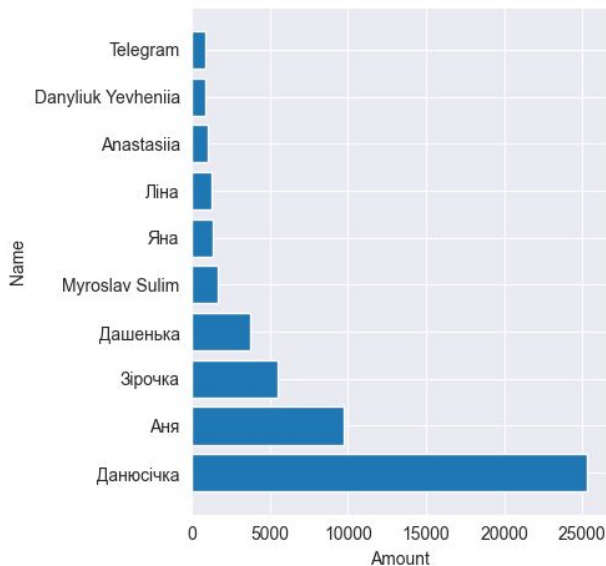
General analysis of my messages

My top-10 emojis

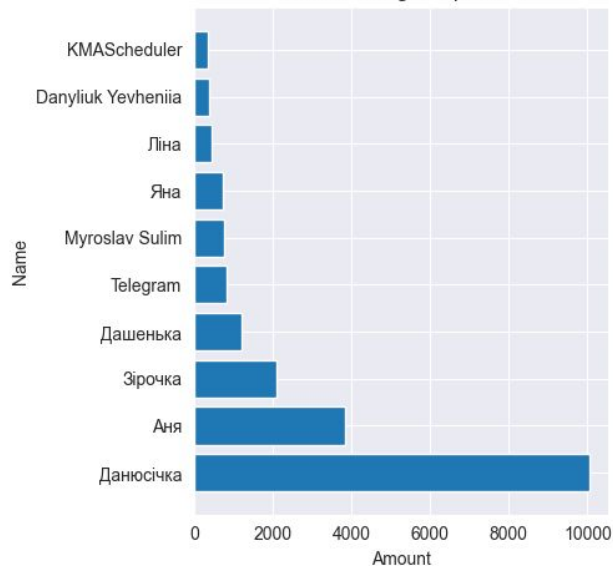


Private dialogs analysis

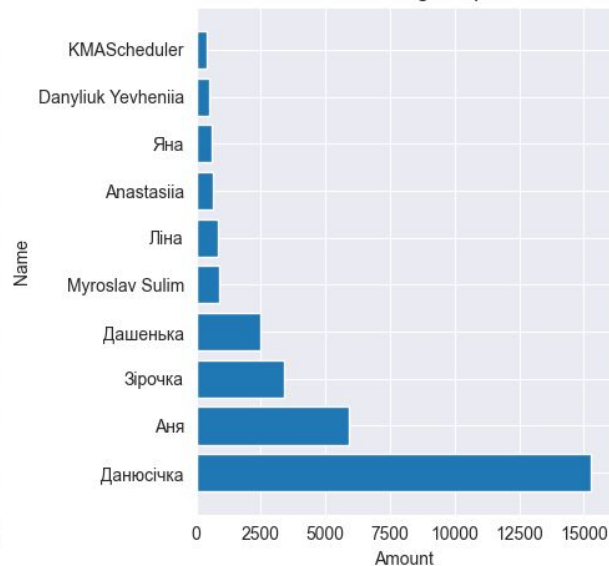
Users with most interections



Send messages top users

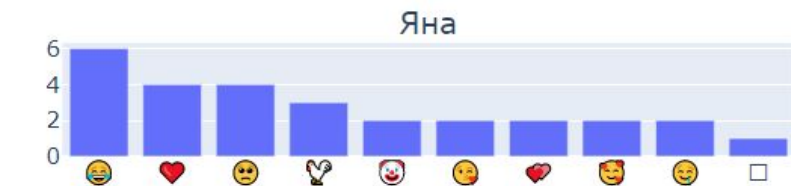
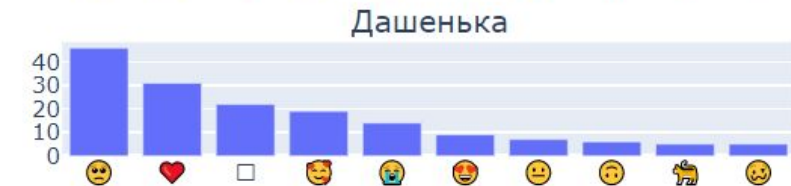
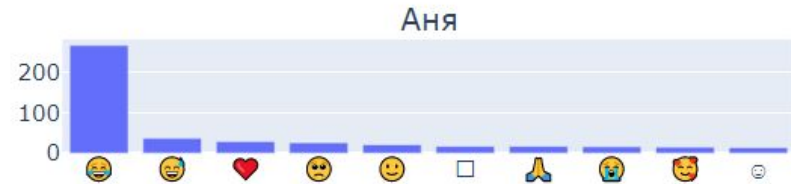
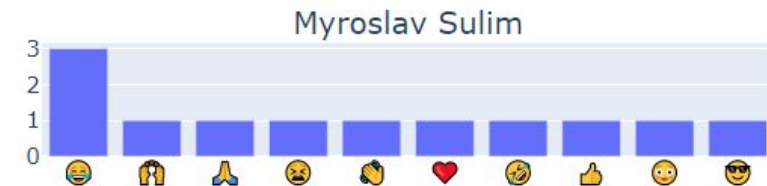
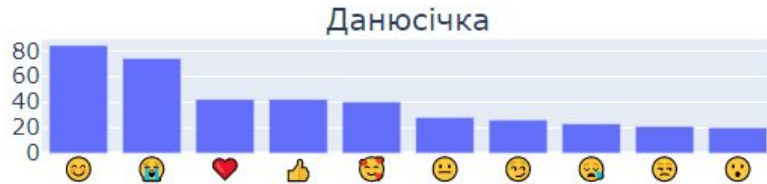


Received messages top users

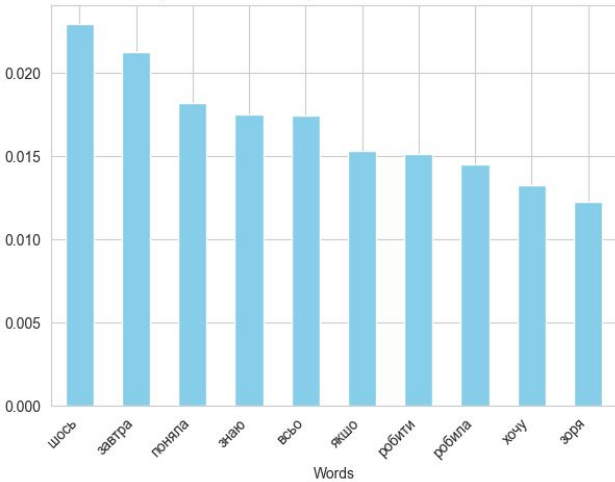


Top emojis of users with whom I interact the most

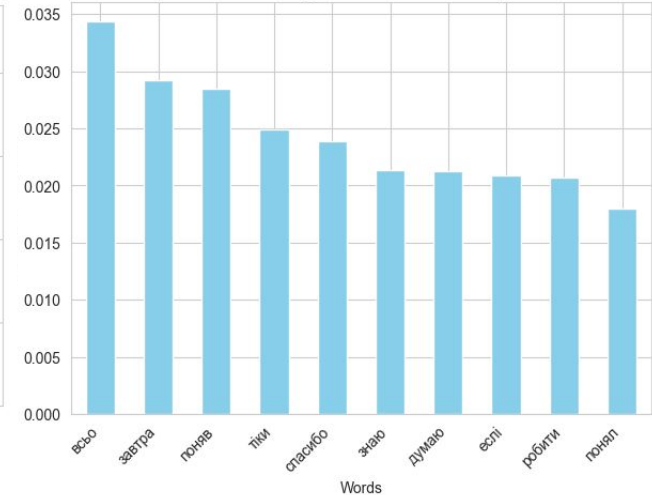
Top-10 Emojis in Received Messages



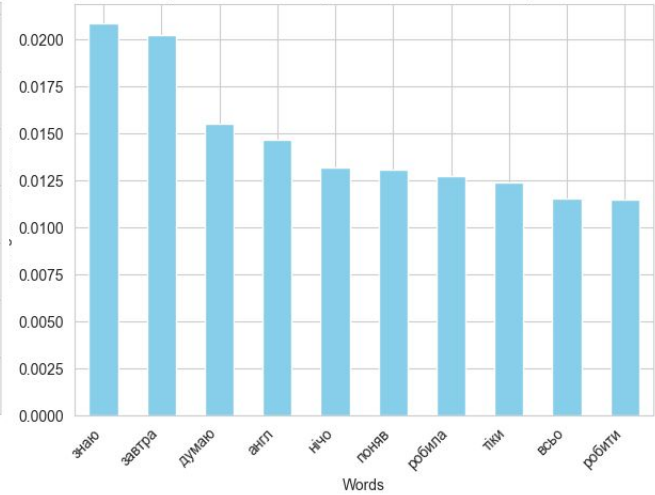
Top 10 Words with Highest TF-IDF Scores for Аня



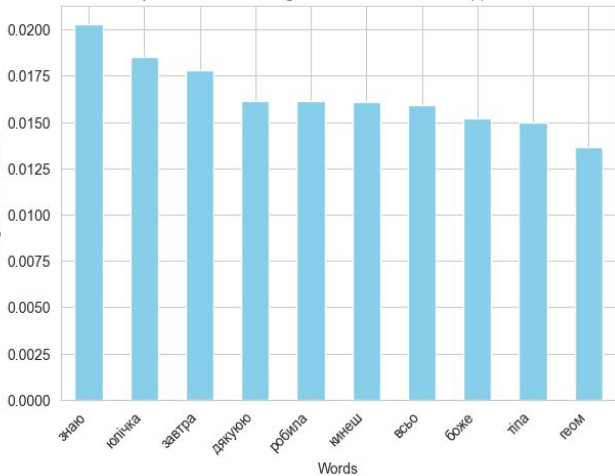
Top 10 Words with Highest TF-IDF Scores for Данюсічка



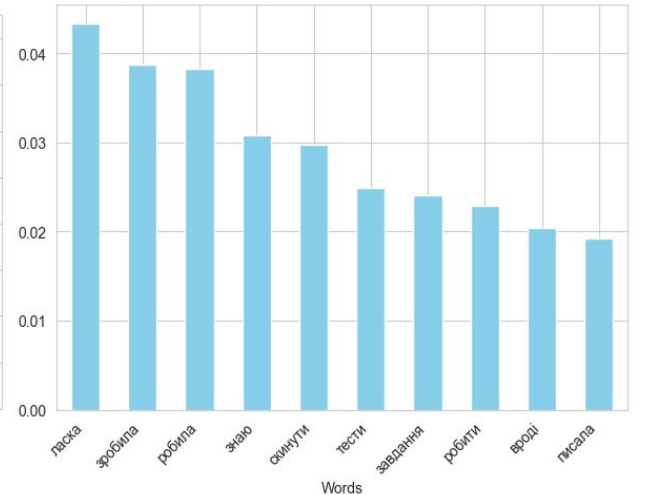
Top 10 Words with Highest TF-IDF Scores for Зірочка



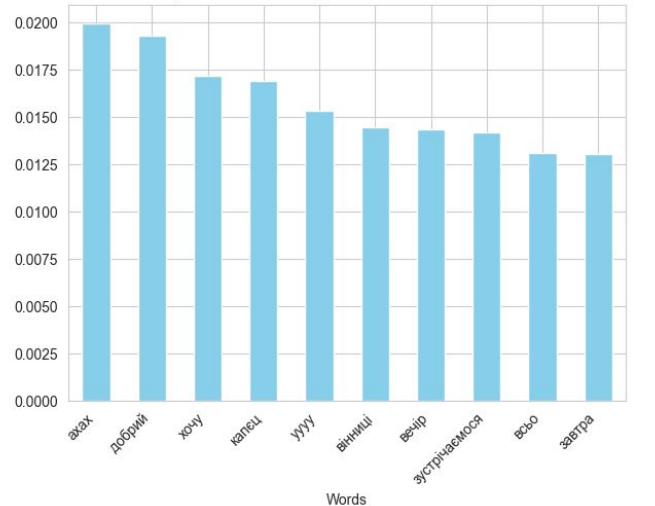
Top 10 Words with Highest TF-IDF Scores for Дашенька



Top 10 Words with Highest TF-IDF Scores for Myroslav Sulim



Top 10 Words with Highest TF-IDF Scores for Яна



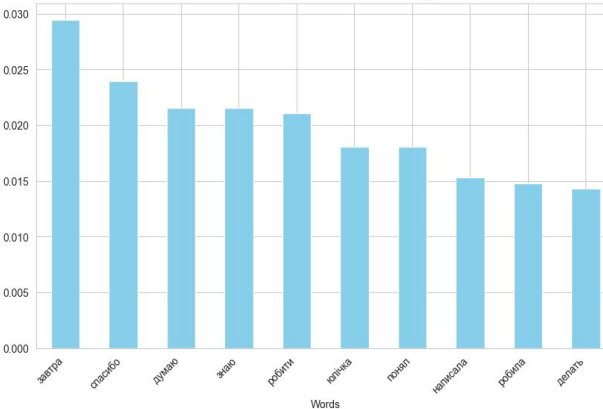
Data cleaning problem

Histograms of words with top tf-idf show us that there are still a lot of words that don't give meaningful information about dialogs.

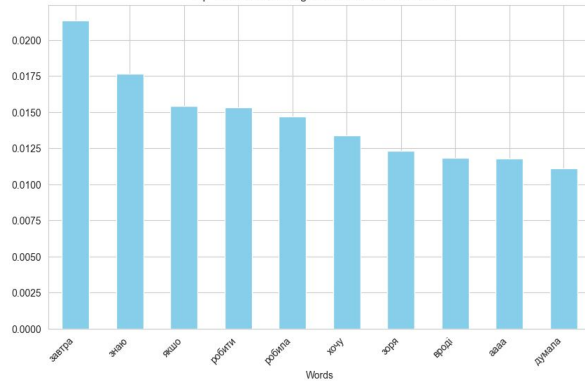
To fix this problem I created new additional dictionary of stop-words.

Updated tf-idf

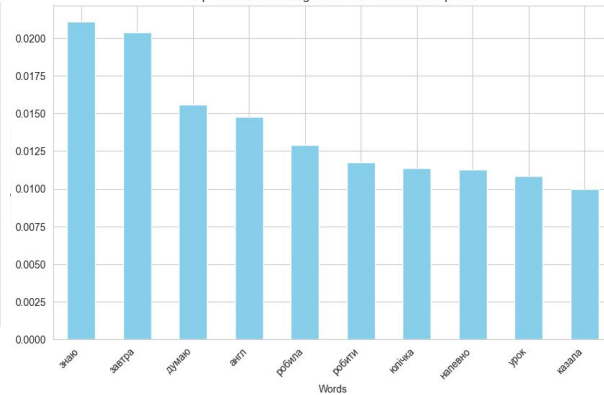
Top 10 Words with Highest TF-IDF Scores for Даниюсінка



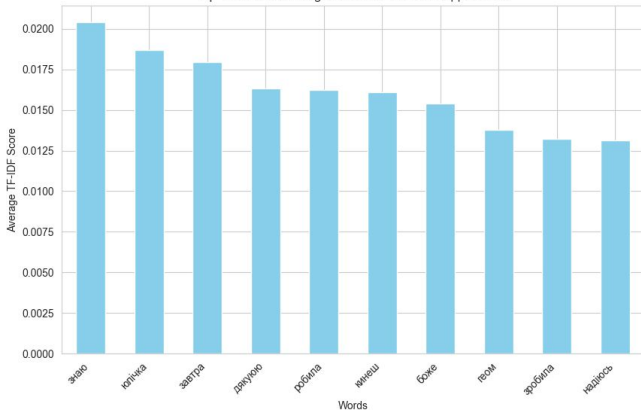
Top 10 Words with Highest TF-IDF Scores for Аня



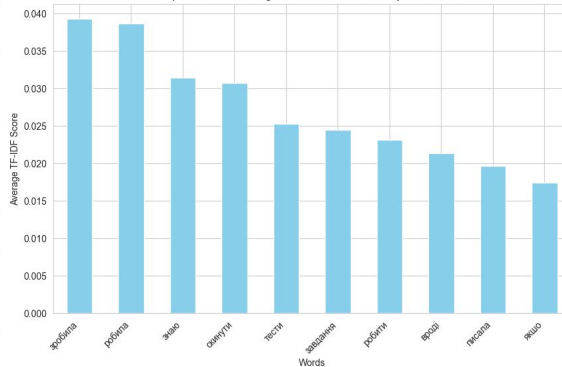
Top 10 Words with Highest TF-IDF Scores for Зірочка



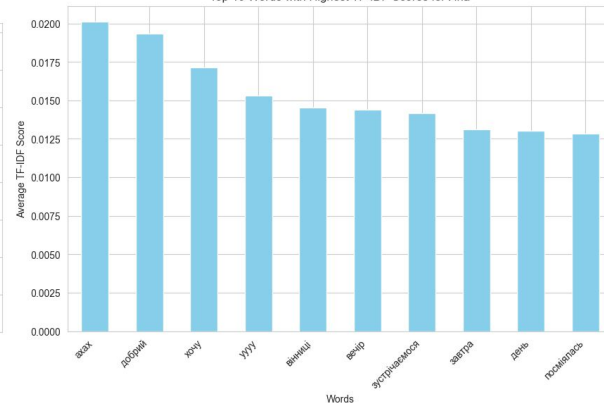
Top 10 Words with Highest TF-IDF Scores for Дашенька



Top 10 Words with Highest TF-IDF Scores for Myroslav Sulim

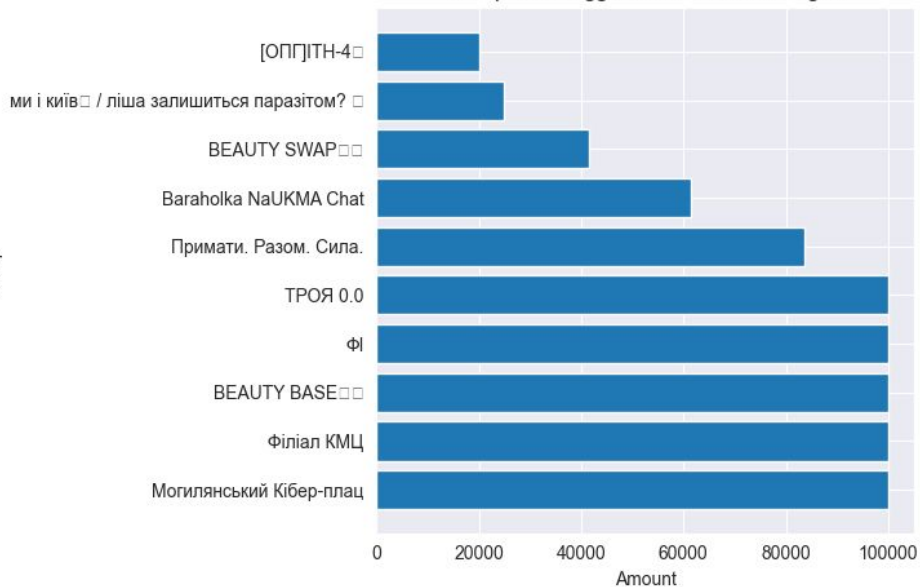


Top 10 Words with Highest TF-IDF Scores for Яна

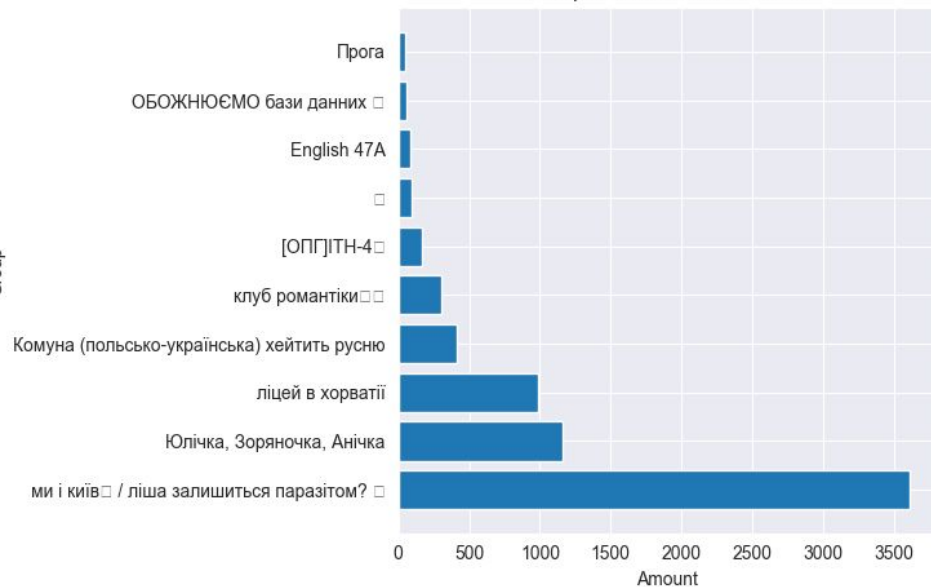


Groups analysis

Groups with biggest amount of messages



Groups where I'm most active



Analysis of Могилянський Кібер-плац

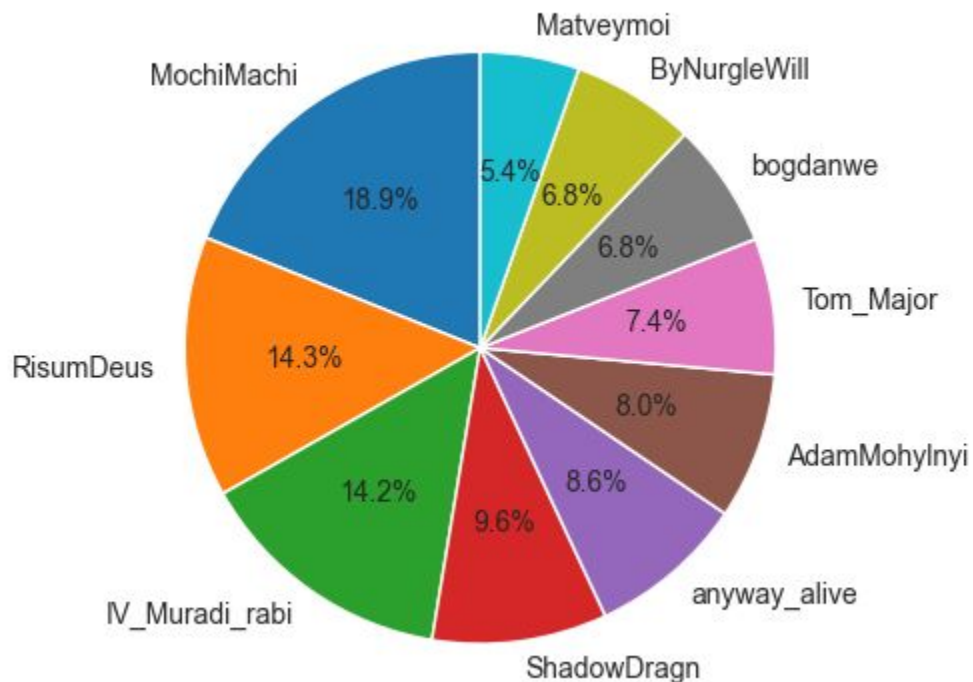
General:

Amount of messages: 100 000

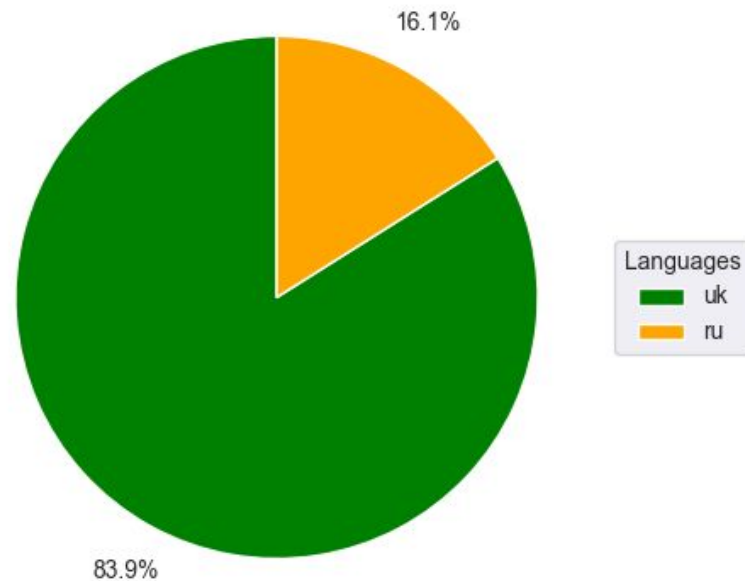
Amount of users: 1613

Analysis of Могилянський Кібер-плац

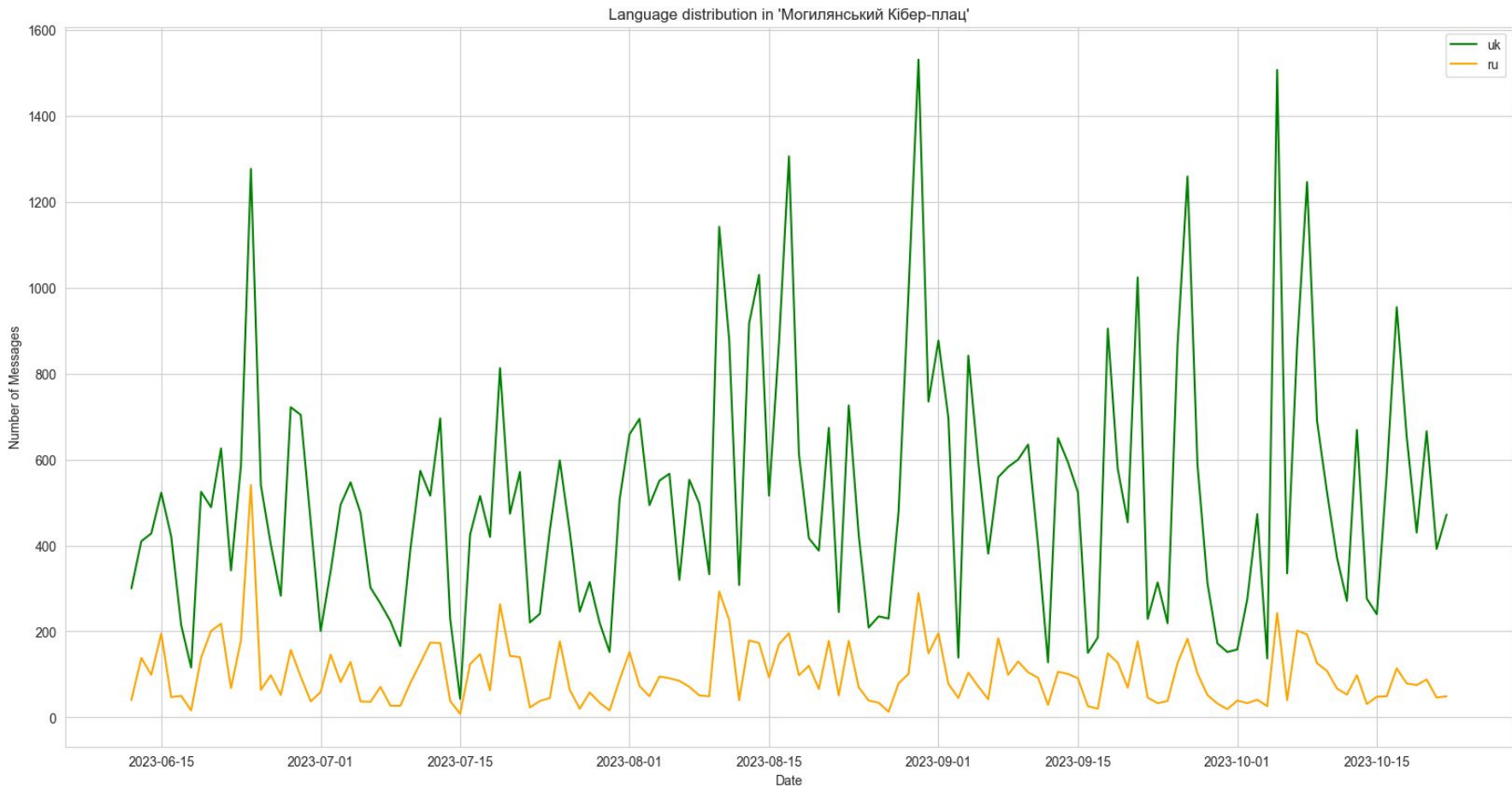
Most active users in 'Могилянський Кібер-плац'



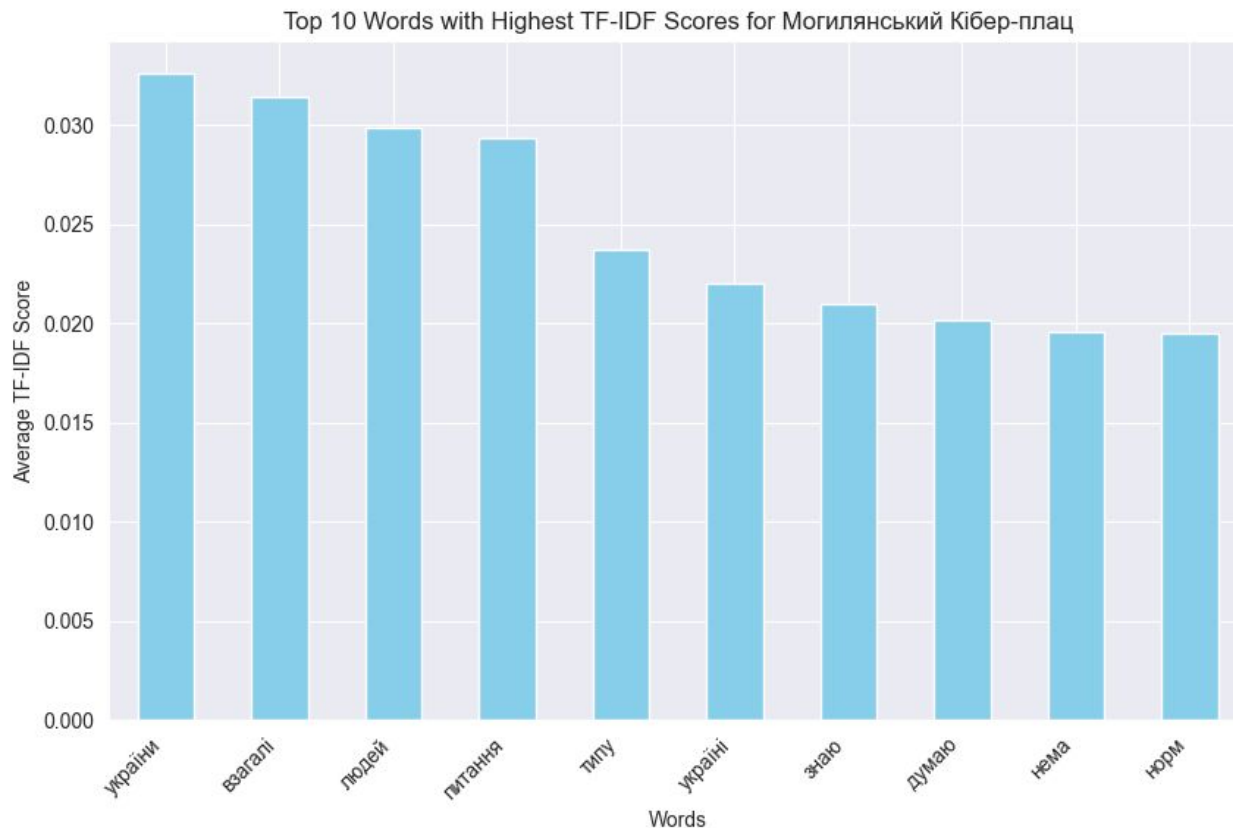
Language distribution in 'Могилянський Кібер-плац'



Analysis of Могилянський Кібер-плац

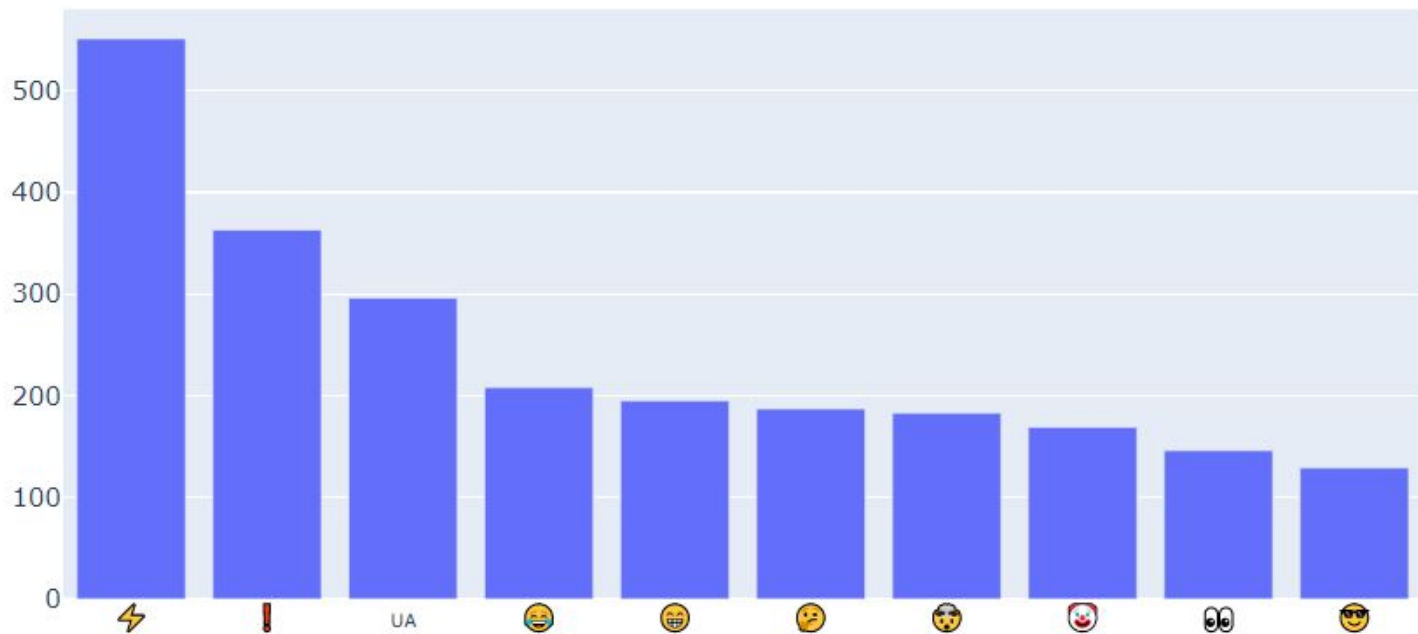


Analysis of Могилянський Кібер-плац



Analysis of Могилянський Кібер-плац

Top-10 emojis of 'Могилянський Кібер-плац'



Analysis of BEAUTY BASE

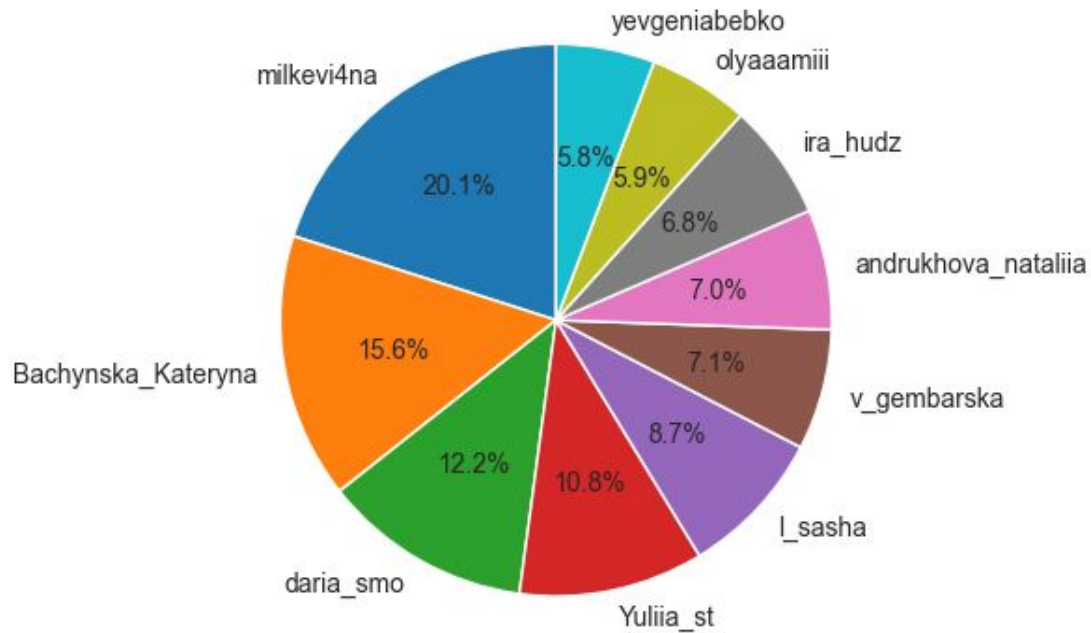
General:

Amount of messages: 100 000

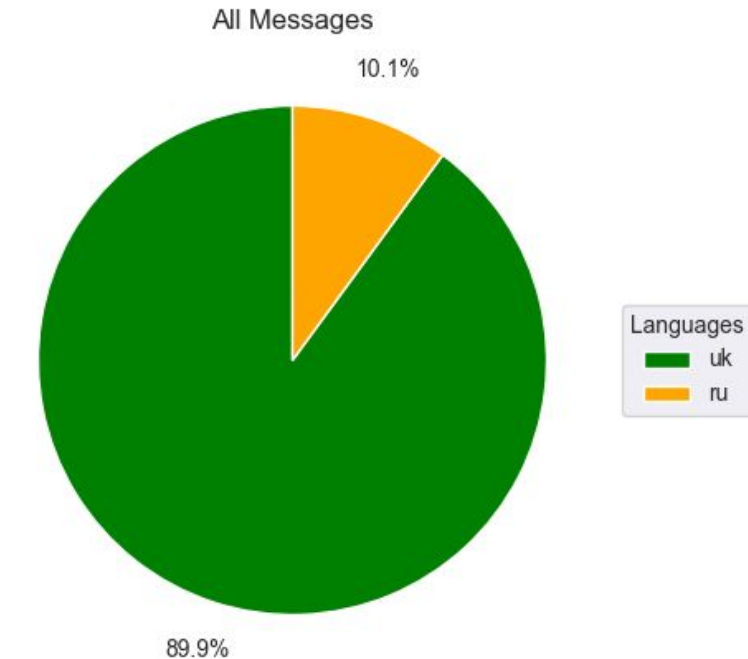
Amount of users: 1613

Analysis of BEAUTY BASE

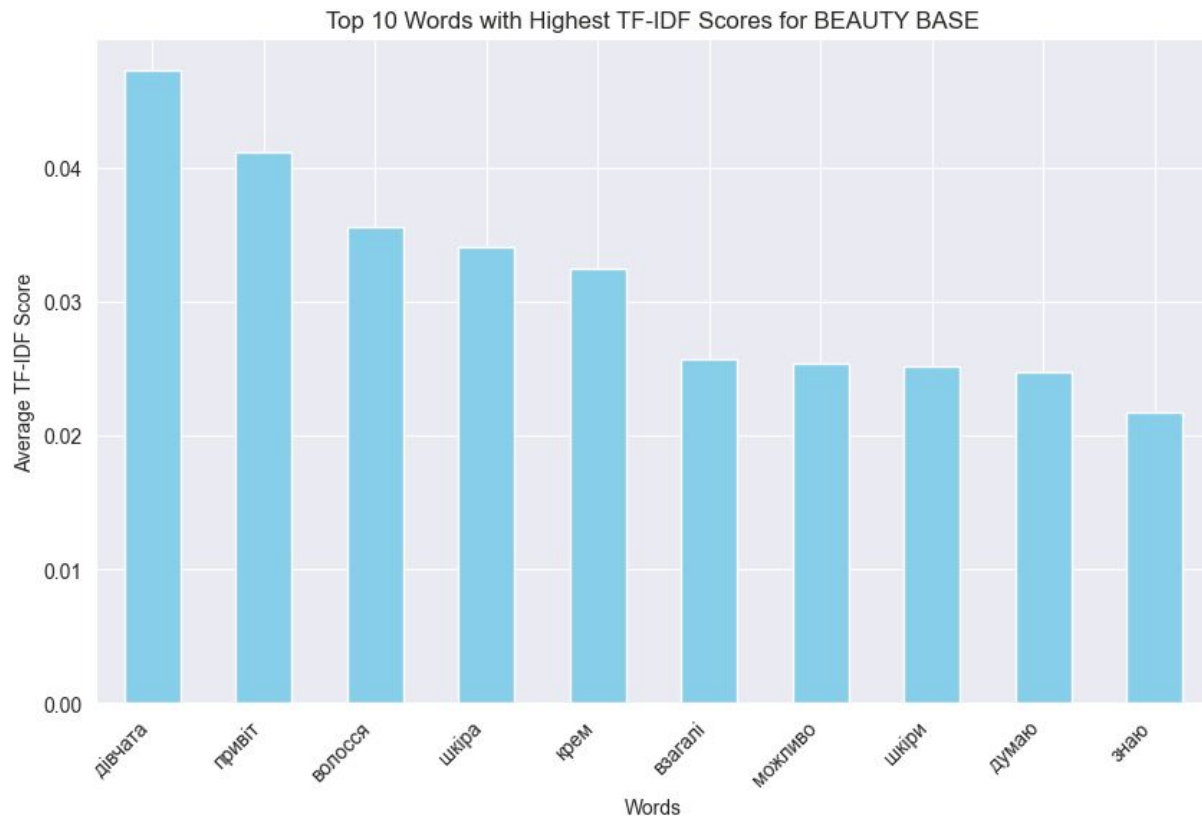
Most active users in 'BEAUTY BASE'



Language distribution in 'BEAUTY BASE'



Analysis of BEAUTY BASE



Analysis of BEAUTY BASE 🧴 📝

Top-10 emojis of 'BEAUTY BASE'



Further work

Improving language detection of dialogs, where messages

Git Repository

[GitLink](#)