



Data Glacier

Your Deep Learning Partner

Exploratory Data Analysis

Hate Speech Detection using Transformers
Deep Learning

https://github.com/kkalyagina/NLP_Project

Agenda

Executive Summary

Problem Statement

Business understanding

Approach

Data understanding

EDA

Recommendations for the model



Data Glacier

Your Deep Learning Partner



Data Glacier

Your Deep Learning Partner

Group Name: individually

Name: Kristina Kaliagina

Email: kalyagina.kristina@gmail.com

Country: Russia

College/Company: Graduated from the University "Higher School of Economics"

Specialization: NLP

Problem description

The task of the project is to classify tweets from Twitter, namely, it is necessary to create a model that will help determine whether a particular tweet belongs to such a type of speech as hate speech. Hate Speech, in simple terms, is offensive language directed at individuals or groups based on their affiliation, interests, and characteristics, such as their religion, nationality, race, color, origin, gender, or other identity factor.

The task is quite difficult because of the inherent complexity of natural language constructs - different forms of hatred, different types of goals, different ways of representing the same meaning.

Business understanding



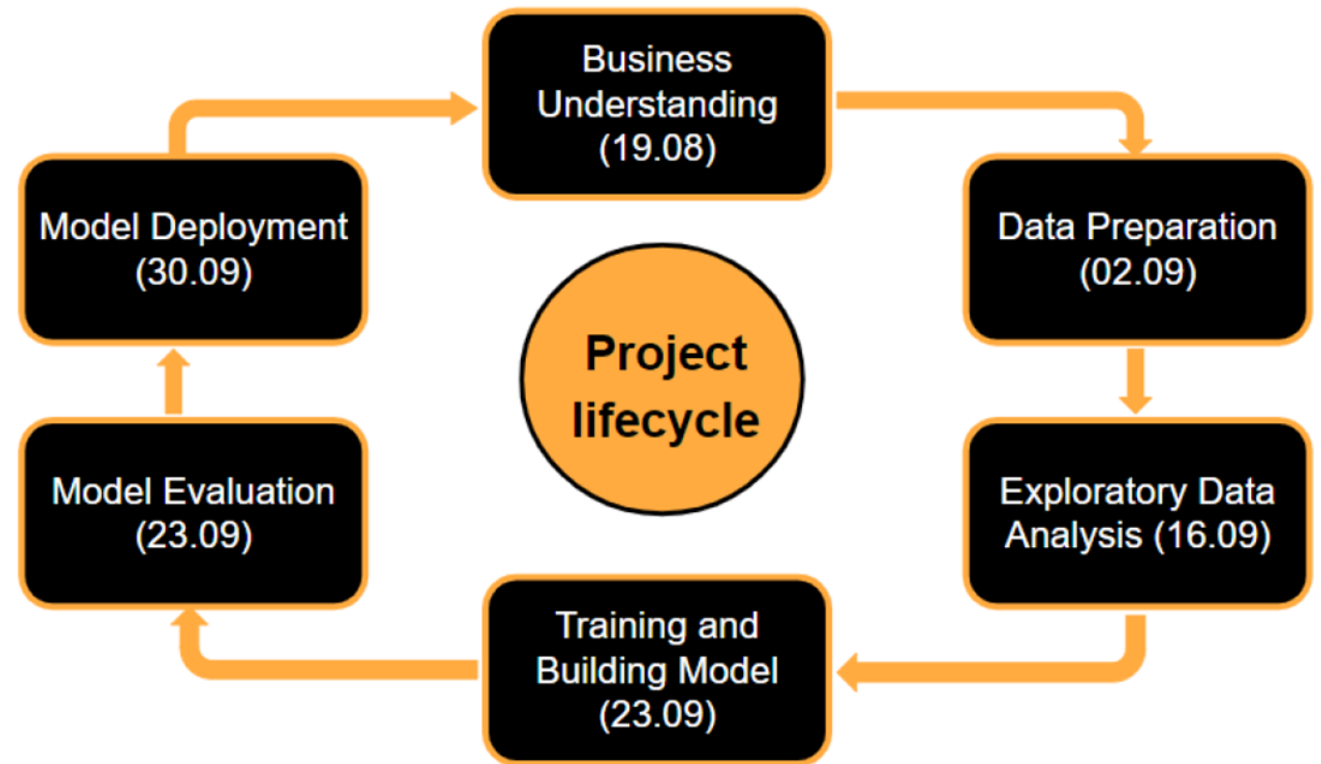
Data Glacier

Your Deep Learning Partner

Social media has become the main driver of social change in global society. The consequences of events taking place in one corner of the world are reflected around the globe in different regions. This is because the vast amount of data generated on these platforms reaches the far corners of the world in the blink of an eye. The developers of these platforms face numerous challenges to make cyberspace as inclusive and healthy as possible. However, in recent years, the phenomena of offensive speech and hate speech have been spreading with greater force. Despite manual efforts, the scale of this problem is so huge that it cannot be solved with coordinated teams. In fact, an automated technique needs to be developed that detects and removes offensive and hateful comments before their harmful effects materialize.

The detection of such hate speech is important for the analysis of public sentiment. User groups in relation to another group, as well as to prevent illegal actions. It's also useful to filter tweets before content recommendations or explore AI chatbots on tweets.

Approach



Data Glacier

Your Deep Learning Partner

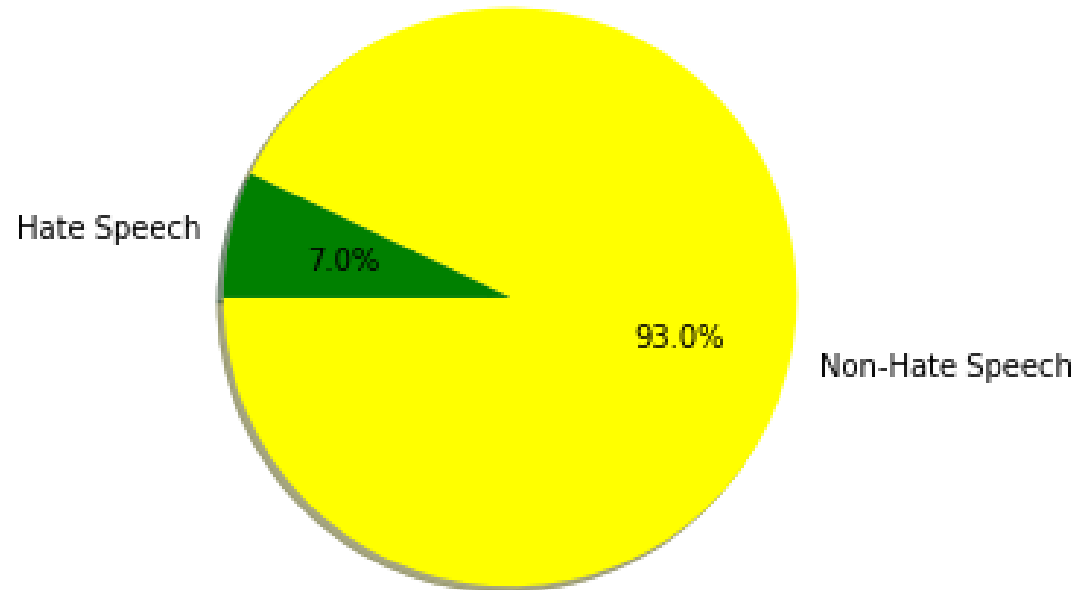
Data understanding

To solve the problem of determining hate speech, data was taken in the form of tweets from the well-known Twitter platform. The data consists of a training and test set. In total, these are 49159 observations, 31961 and 17197, respectively.

Data problems:

- Unbalanced classes
- Duplicate retweets
- Passes, extra spaces
- Lots of punctuation
- @usernames and #hashtags
- Slang
- Interjections
- Greek symbols

Share of hate speech tweets

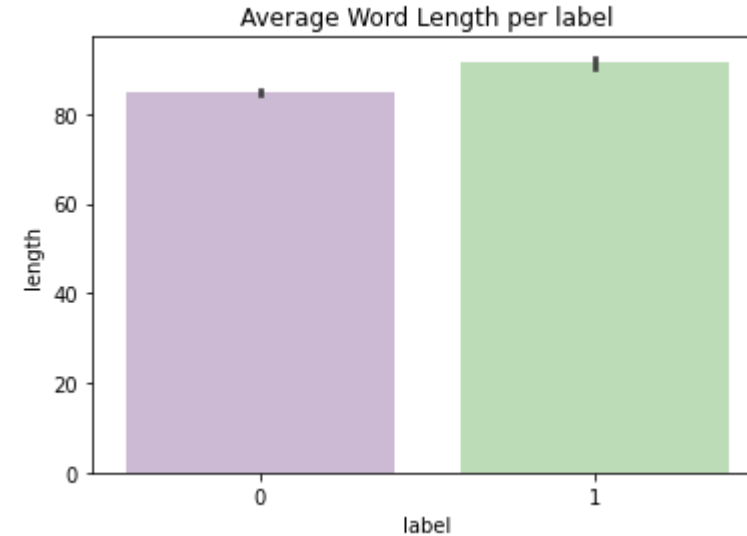
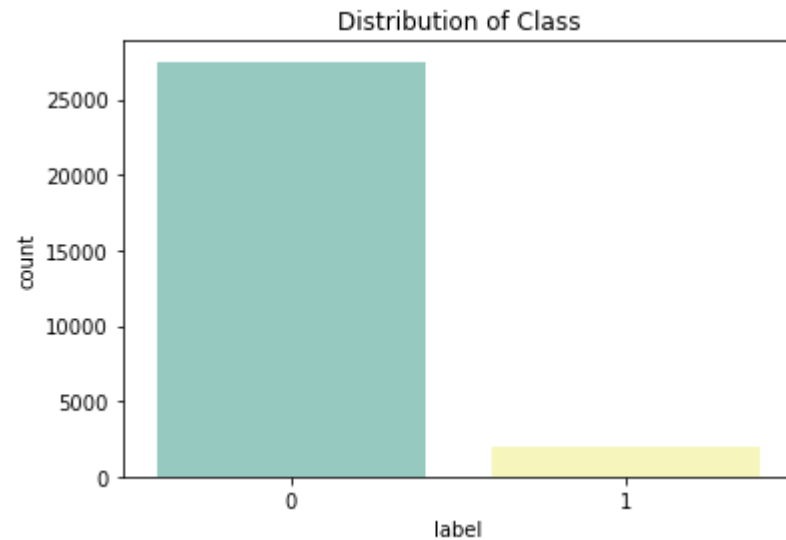


Data Glacier

Your Deep Learning Partner

Exploratory Data Analysis

There are more non hate speech than hate speech tweets. The level of hate speech tweets is 7% of the total number of observations.



The average length of words in tweets containing hate speech is slightly longer than in regular tweets.

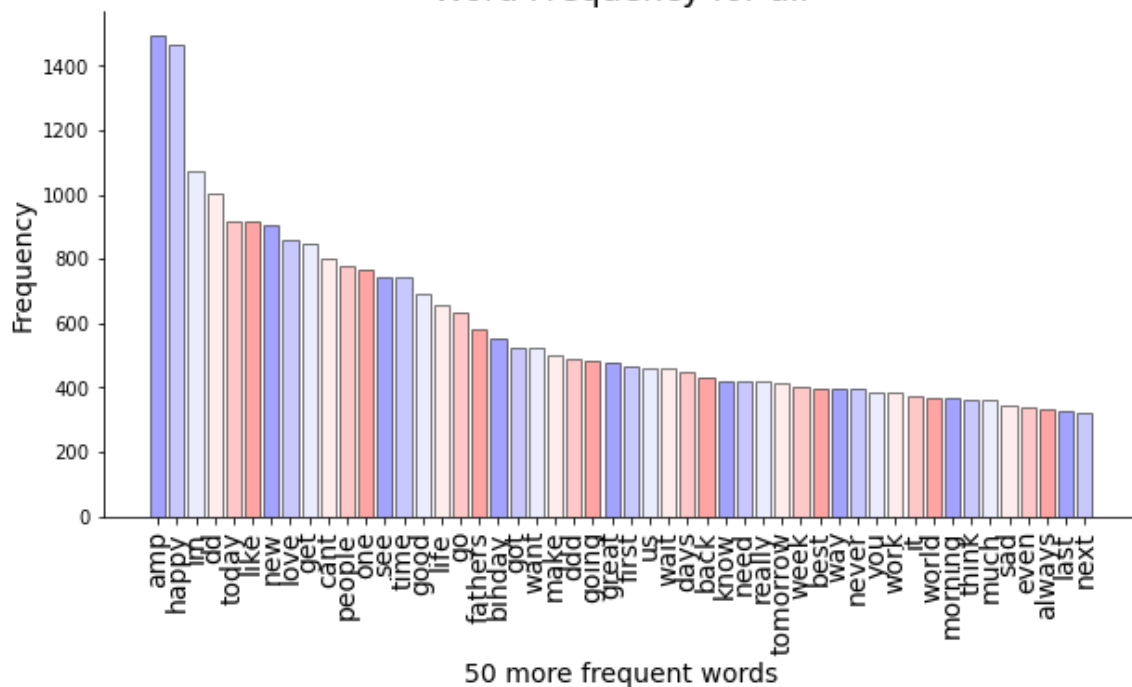


Data Glacier

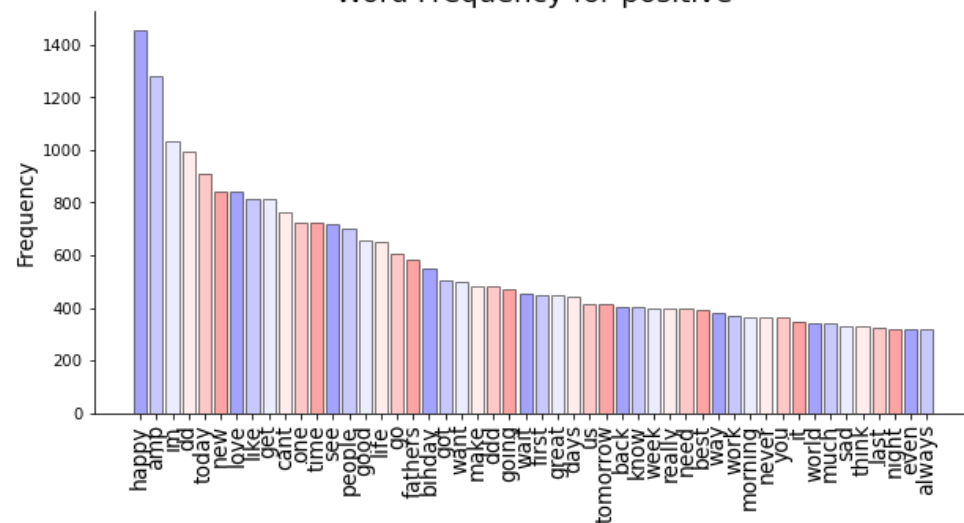
Your Deep Learning Partner

Word frequency

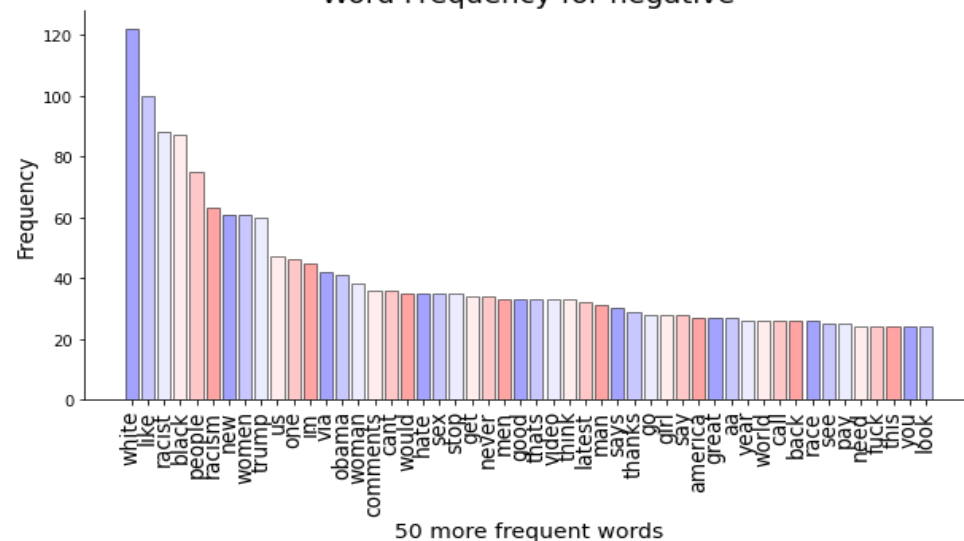
Word Frequency for all



Word Frequency for positive



Word Frequency for negative

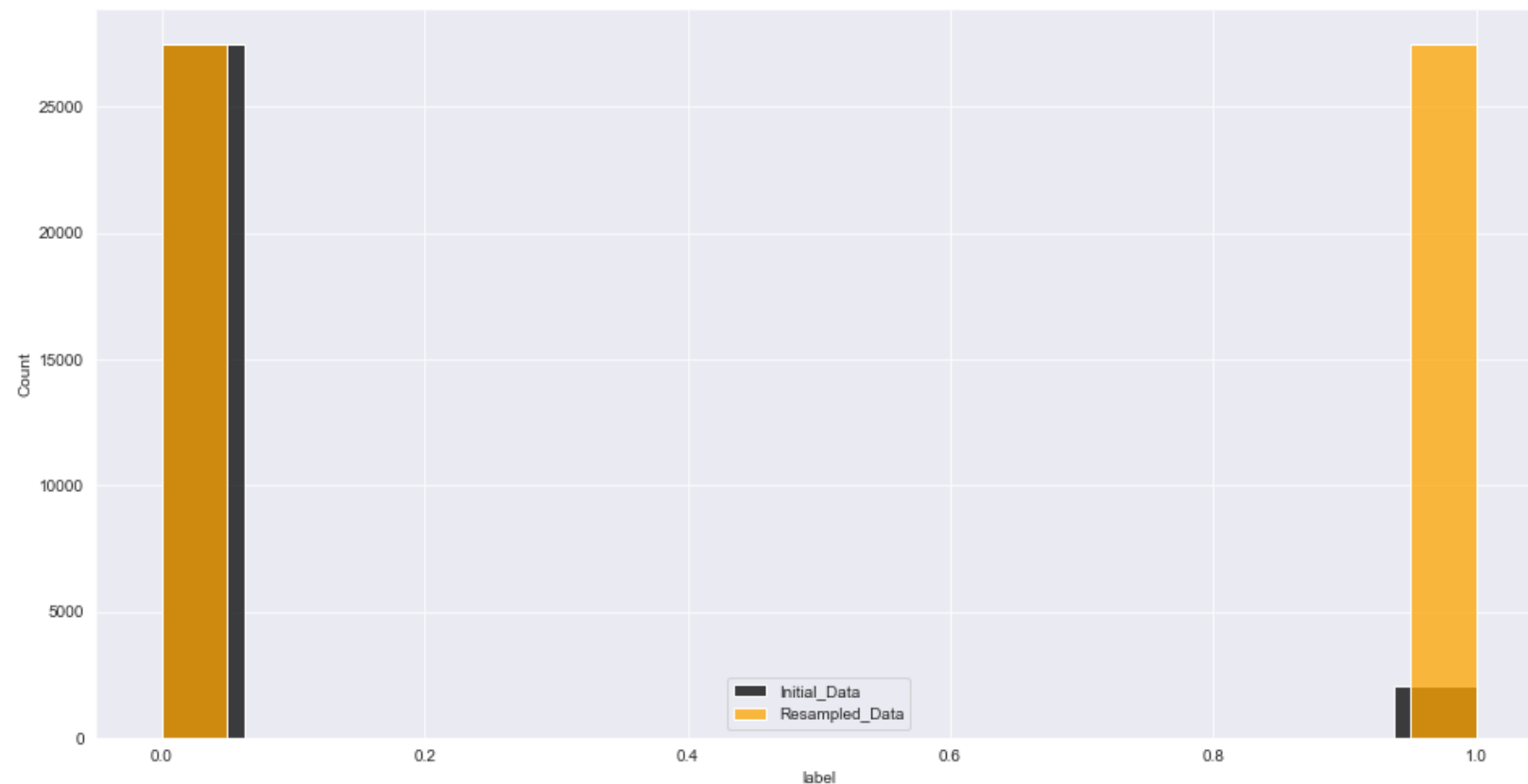


Data Glacier

Your Deep Learning Partner

Unbalanced classes

There is an imbalance of classes in the dataset, so we managed to restore the balance with the help of resampling. Since when building a model, such a ratio can spoil our results.

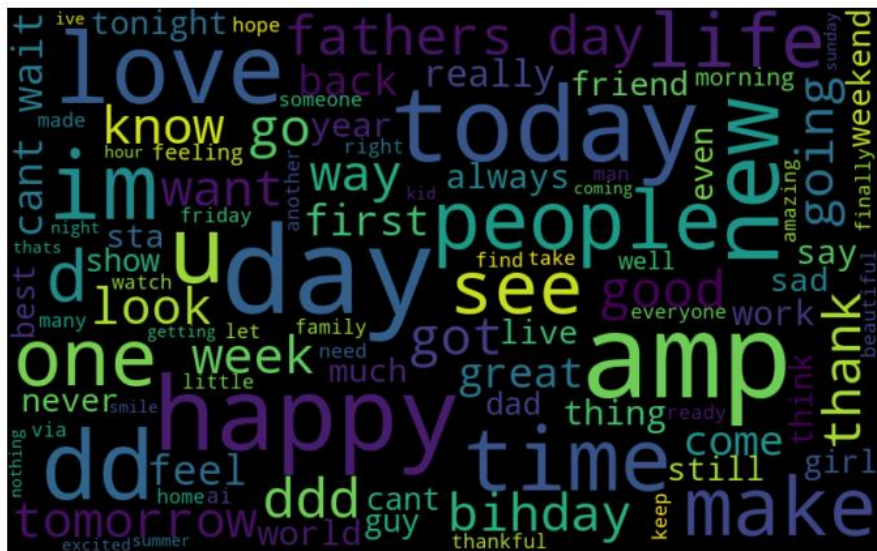


Data Glacier

Your Deep Learning Partner

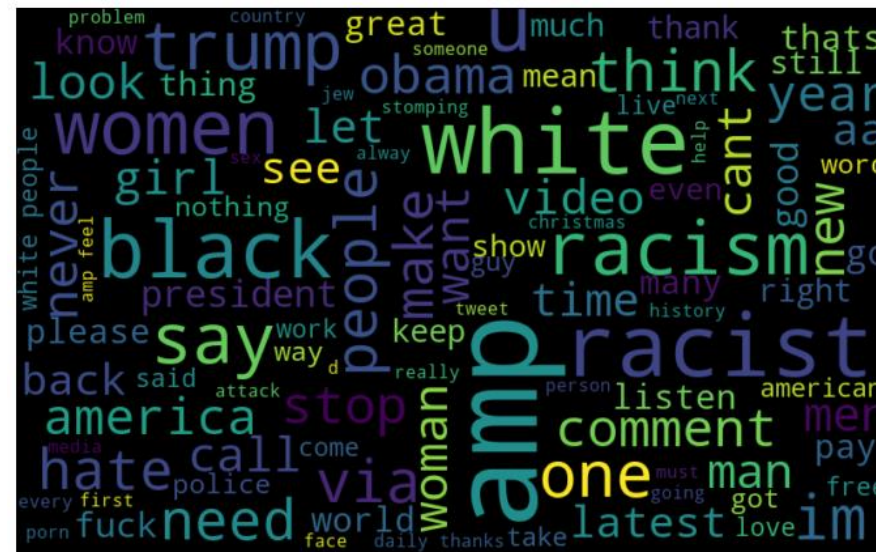
Wordcloud for tweets before resampling

Non-hate speech



We can see that word's common in positive comments are: love, happy, friend, life, today, day, thank, time, see, new, people, one, i'm, fathers day, good and so on

Hate speech



We can see that word's common in hate comments are: trump, hate, white, black, racist, racism, allahsoil, obama, women, never, america, stop and so on

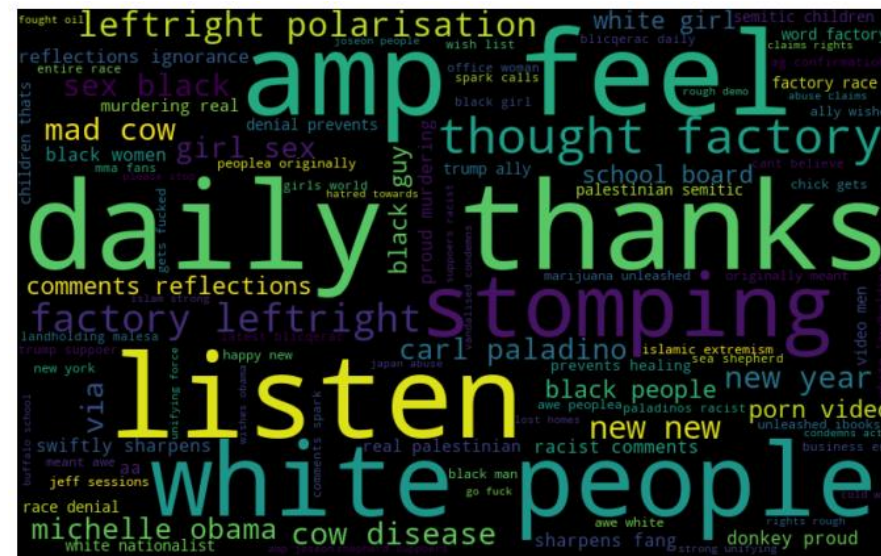
Data Glacier

Wordcloud for tweets after resampling

Non-hate speech



Hate speech



There are almost the same positive words as in the situation before the resample

The last image observes words and people related to religion, politics and sex



Your Deep Learning Partner

Recommendations for the model

Tokenizer

Embedding

LSTM

Loss: binary_crossentropy

Optimizer: Adam

Metrics: Accuracy

Transformer Block

Keras Model



Data Glacier

Your Deep Learning Partner



Data Glacier

Your Deep Learning Partner

Thank you for attention!