# NLP Project phase 1

Saba Rayehi - Sohrab Namazi Nia
Iran University of Science and Technology

# Toxic comment definition

The goal of this project is to classify comments into two groups, toxic comments, and safe comments. A comment must be classified as toxic if it is rude, disrespectful, or unreasonable.

# Toxic comment Dataset

We searched for many different datasets on the internet and finally chose one of them that was a better generalization of real-world comments.

The raw dataset can be found here:

[Link](#)

# Dataset Structure

The raw dataset is written into CSV files. Each element of the training dataset, consists of a clear label, stating whether the data is toxic or not. Also, there are 5 parameters for each element that are mentioned here:
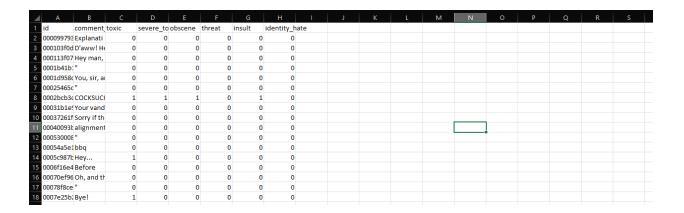
**Severe_toxic**

**Obscene**

**Threat**

**Insult**

**Identity_hate**

The values for the mentioned parameters are binary, meaning that for example, a comment is whether a threat or not. Also, each comment has a unique ID.

You can see a demo of what the CSV file for training data looks like in the following figure:

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | id | comment | toxic | severe_to | obscene | threat | insult | identity_hate | | | | | | | | | | | |
| 2 | 000099793 | Explanati | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 3 | 000103f0d | D'aww! H | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 4 | 000113f07 | Hey man, | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 5 | 0001b41b: | " | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 6 | 0001d958c | You, sir, a | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 7 | 00025465c | " | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 8 | 0002bcb3c | COCKSUCI | 1 | 1 | 1 | 0 | 1 | 0 | | | | | | | | | | | |
| 9 | 00031b1e! | Your vand | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 10 | 00037261f | Sorry if th | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 11 | 00040093& | alignment | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 12 | 00053000& | " | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 13 | 00054a5e1 | bbq | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 14 | 0005c987k | Hey... | 1 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 15 | 0006f16e4 | Before | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 16 | 00070ef96 | Oh, and th | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 17 | 00078f8ce | " | 0 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |
| 18 | 0007e25b: | Bye! | 1 | 0 | 0 | 0 | 0 | 0 | | | | | | | | | | | |

## Dataset statistics

We have acquired different asked parameters about the dataset in our code. For example sentence count, word count, etc.

We use nltk and its functions :

For separation sentences —-> we use function sent_tokenize

For separation words—--> we use function word_tokenize

For clean _data:

To clean the data, we first removed the spaces and stopwords then the nonsense words and characters