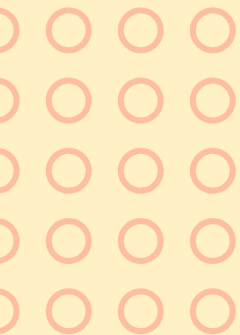




Multilingual Toxic Comment Classification

Author: Drenizë Rama



01

**INTRODUCTION
TO THE PROBLEM**

02

THE DATA

03

NLP MODELS

04

**PRACTICAL
APPLICATION**

TABLE OF CONTENTS






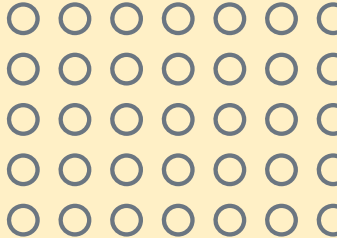
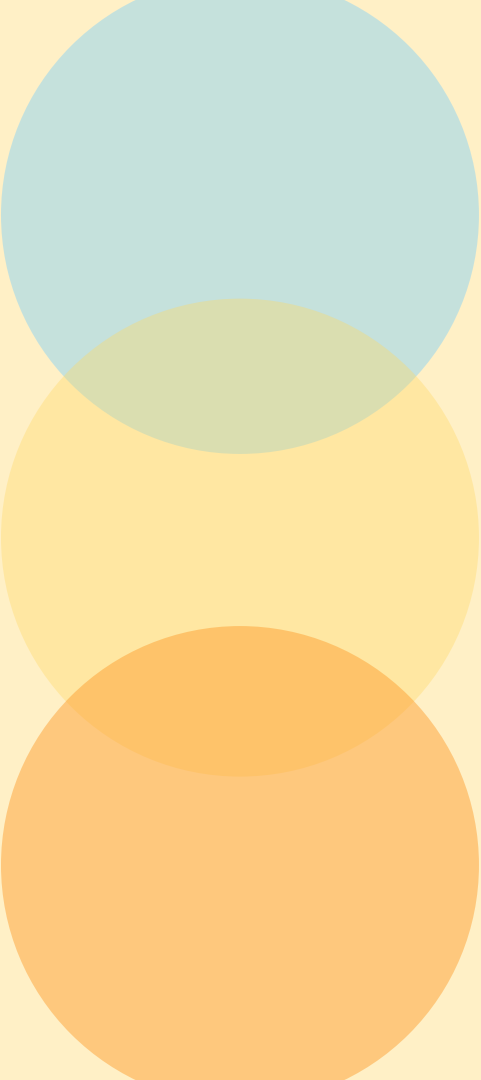
"If you want to feel shit about the state of the world and everyone in it, Facebook comments are a great place to start." a [LinkedIn](#) article states.

How is it possible, though, that the very democratic element of the internet - the **comment section** - has become such a tool of cruelty?

Even the number one internet culture curator [WIRED](#) magazine wipes the slates clean, by stating:

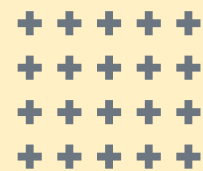
INTERNET RULE #1: Never read the comments.





The goal of this project is to build a natural language algorithm that classifies text input (social media comments) into toxic and non-toxic categories in multiple languages.

In order to achieve this, I decided to go with pretrained RNN models like BERT or LASER, which include a variety of languages.



EDA - The Data

4 datasets:

```
1 train_data.head(2)
```

	id	comment_text	toxic	severe_toxic	obscene	threat	insult	identity_hate
0	0000997932d777bf	Explanation\nWhy the edits made under my usern...	0	0	0	0	0	0
1	000103f0d9cfb60f	D'aww! He matches this background colour I'm s...	0	0	0	0	0	0

```
1 valid_data.head(2)
```

	id	comment_text	lang	toxic
0	0	Este usuario ni siquiera llega al rango de ...	es	0
1	1	Il testo di questa voce pare esser scopiazzato...	it	0

```
1 test_data.head(2)
```

	id	content	lang
0	0	Doctor Who adlı viki başlığına 12. doctor olar...	tr
1	1	Вполне возможно, но я пока не вижу необходимо...	ru



EDA - The Data

And an extended dataset for further exploration:

```
1 data.columns
```

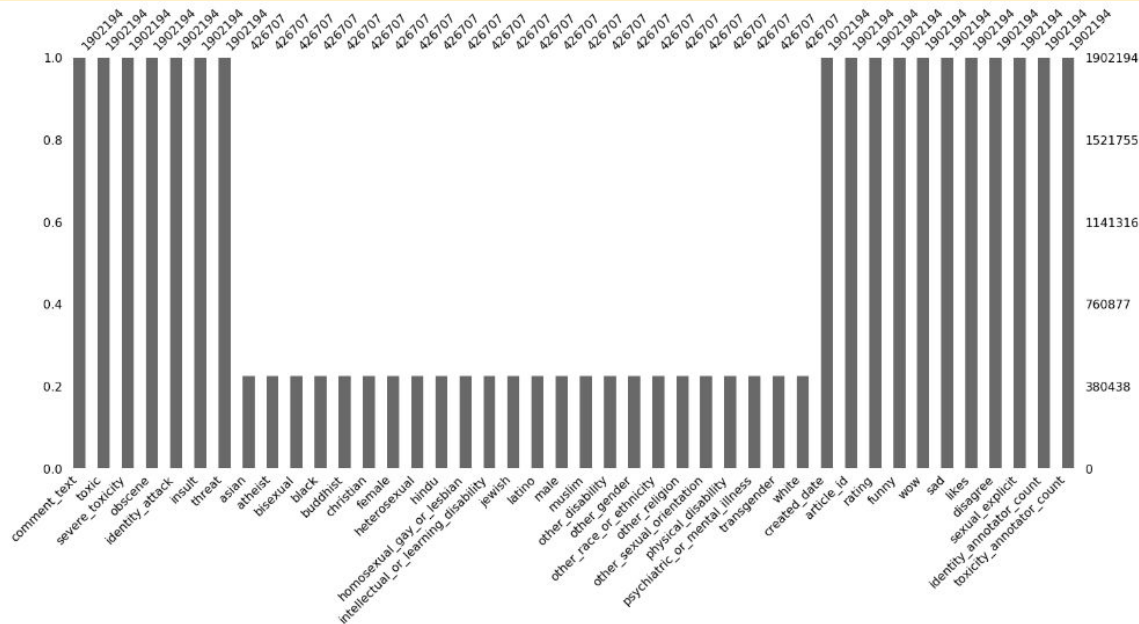
```
Index(['comment_text', 'toxic', 'severe_toxicity', 'obscene',  
      'identity_attack', 'insult', 'threat', 'asian', 'atheist', 'bisexual',  
      'black', 'buddhist', 'christian', 'female', 'heterosexual', 'hindu',  
      'homosexual_gay_or_lesbian', 'intellectual_or_learning_disability',  
      'jewish', 'latino', 'male', 'muslim', 'other_disability',  
      'other_gender', 'other_race_or_ethnicity', 'other_religion',  
      'other_sexual_orientation', 'physical_disability',  
      'psychiatric_or_mental_illness', 'transgender', 'white', 'created_date',  
      'article_id', 'rating', 'funny', 'wow', 'sad', 'likes', 'disagree',  
      'sexual_explicit', 'identity_annotator_count',  
      'toxicity_annotator_count'],  
      dtype='object')
```

```
1 data.shape
```

```
(1902194, 42)
```



The Data

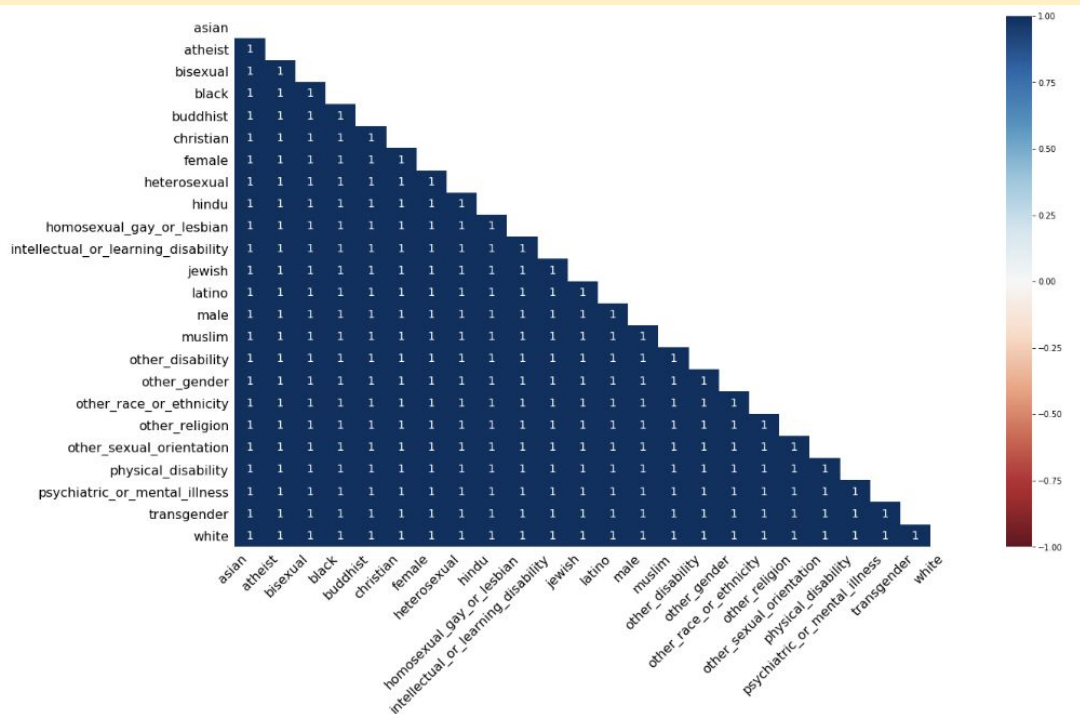


What do we see in the bar chart above?

- There is an equal amount of missing values in all of the identity features.

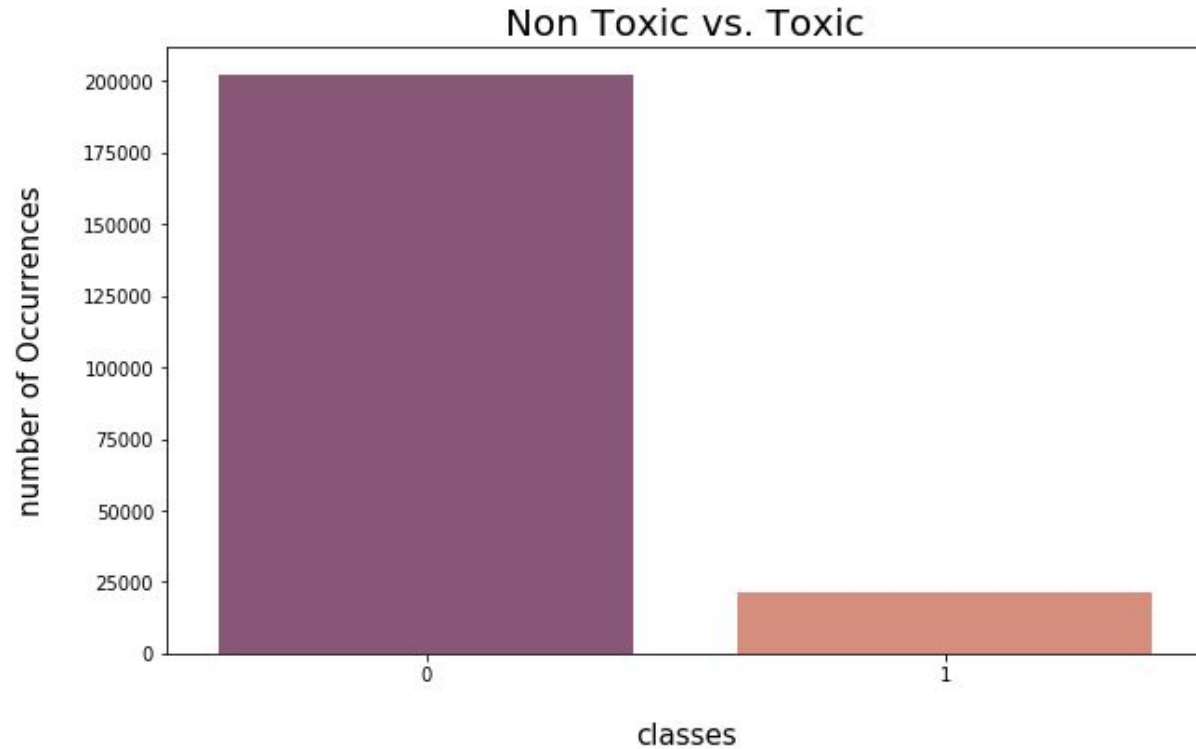
This may indicate that the type of identity hate has been evaluated only later on. However, it is more likely that two datasets have been joined.

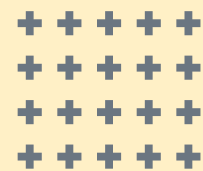
The Data



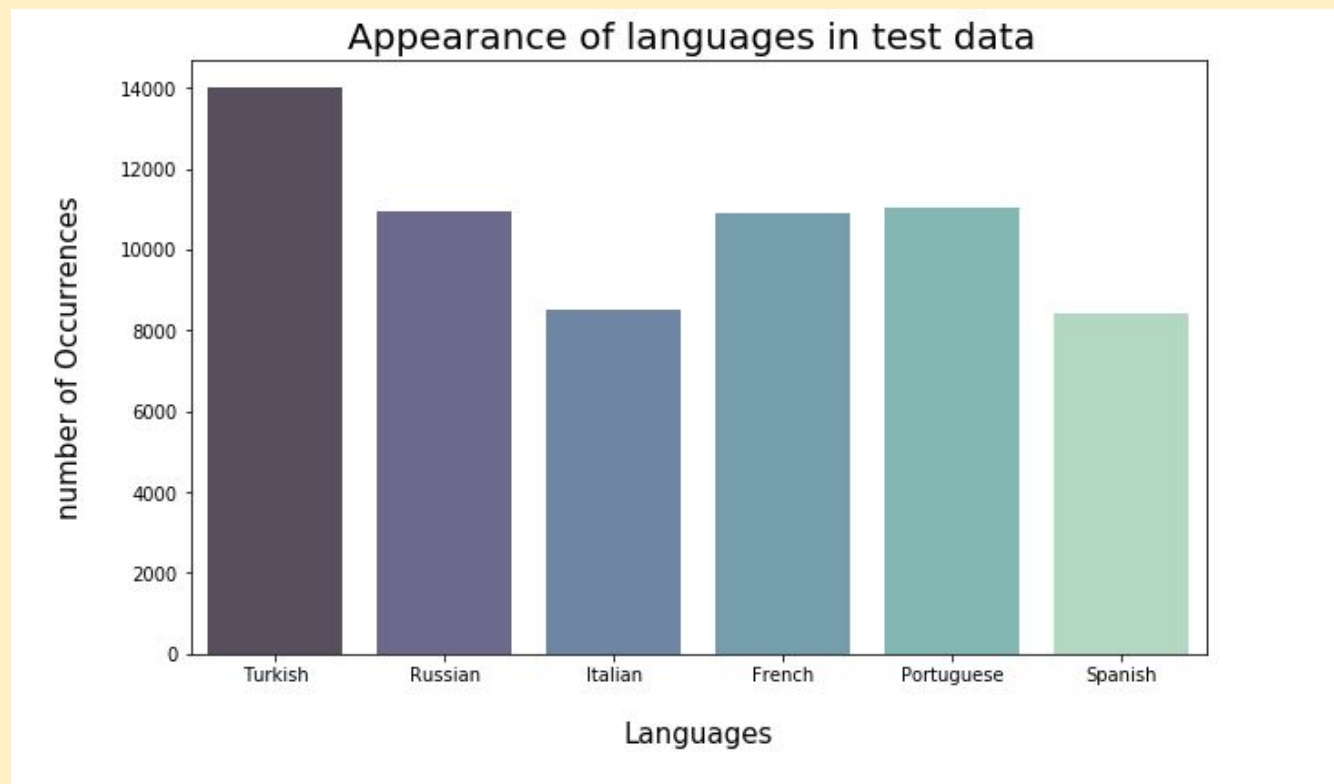
The perfect value of 1 in the *missing values heatmap* confirms that if one identity feature is present the other identity features are also present, suggesting that they have been added at the same time. This is a systematic error, so just going with the 'complete cases analysis' and reducing the sample to complete cases won't do. This ignorance might reduce analysis precision significantly. And that's surely not what we want.

The Data

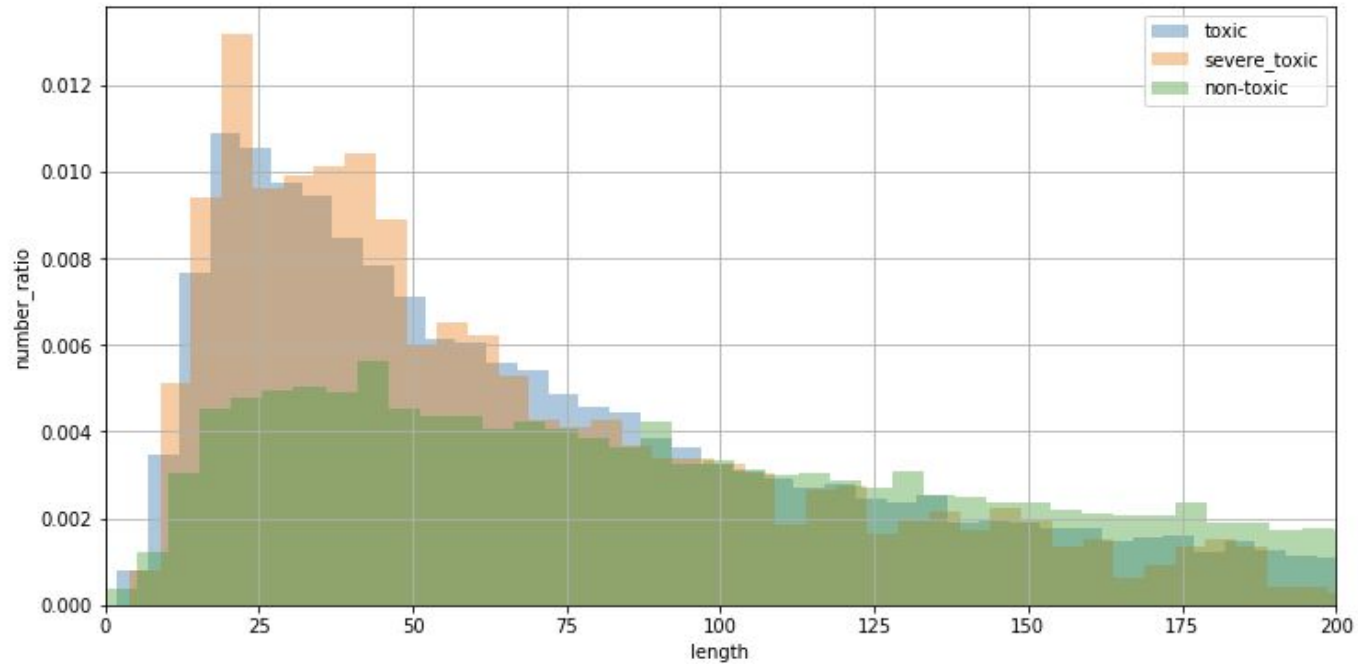




The Data

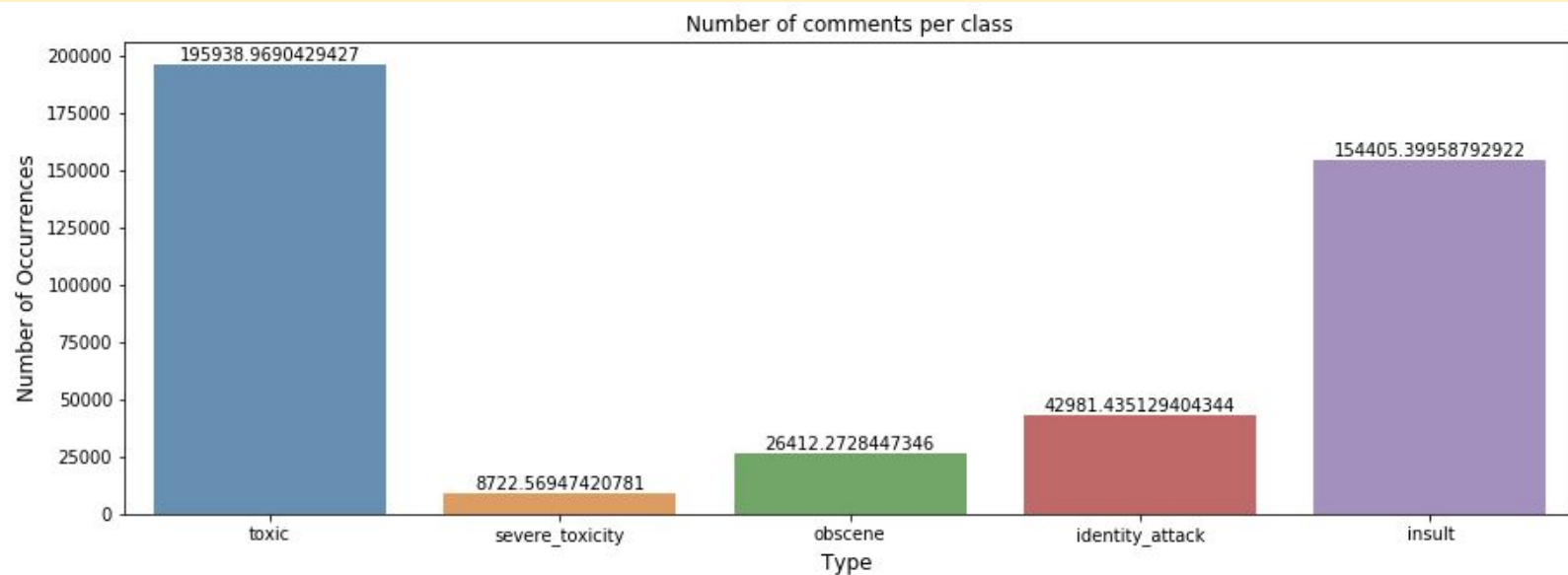
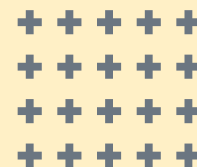


The Data



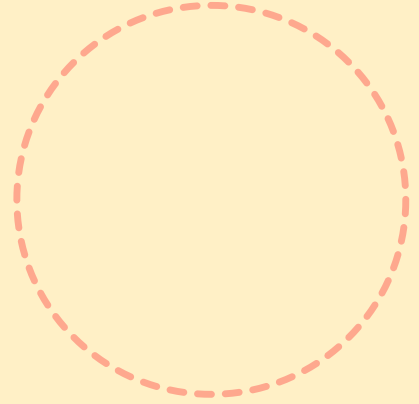
Length of comments by toxicity

The Data





THANKS!



CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon, and infographics & images by Freepik.

