



Data Analytics

Analysis of Toxic Comments on Social Media Platforms And Predict the Level of Toxicity

CAPELLI Guillaume

Nov, 2023

Table of content

Table of content	2
Introduction:.....	3
Use Case Overview:	3
Project Goal:	3
Objective:	3
Data and data sources	4
YouTube.....	4
Reddit:.....	4
Kaggle:	4
BigQuery:	4
Data collection	5
GDPR	6
Data Cleaning and Exploratory Data Analysis	7
Data Cleaning:	7
EDA:	7
Here's some visualisation examples:	8
Data base type selection	13
Entities. ERD	14
Database's Queries.....	15
Database's Queries – Outputs	16
Exposing Data via API	17
API – Homepage.....	18
References	19

Introduction:

Understanding Online Toxicity

In a world where online interactions are increasingly shaping our perceptions and everyday experiences, the presence of toxic comments on social media has emerged as a critical concern. For me, this analysis is not just a scholarly pursuit but a necessary step towards fostering a safe digital environment. Toxic comments can deter meaningful discourse, perpetuate negativity, and even cause significant psychological harm. By analyzing these comments, I aim to unveil the patterns and triggers of online hostility, providing insights that can help platforms and communities curtail this pervasive issue.

Use Case Overview:

The digital landscape offers an unparalleled platform for free expression. Yet, the anonymity and detachment provided by screens often lead to an increase in hostile and aggressive communication. My focus is on identifying and understanding these patterns of toxicity to mitigate their impact and support the creation of more respectful online communities.

Project Goal:

The primary goal of my project is to dissect the complex dynamics of toxic comments across various social media platforms. The end objective is to collate these findings to inform and train a machine learning model that can detect and categorize toxicity autonomously, thereby aiding platform moderators and community managers.

Objective:

I have set specific objectives to guide my analysis:

- .To analyze trends in comments over time and assess how user engagement with content evolves.
- .To investigate the relationship between comment attributes and user interaction, such as likes and replies.
- .To categorize comments into different levels of toxicity, which will serve as foundational data for training my predictive model.

Data and data sources

YouTube: Utilizing Selenium, I scraped comments from select videos that were known for polarizing content. I chose YouTube because it's a widely used platform where video content often elicits strong opinions and reactions, which are reflected in the comments section. This dataset offers a window into user reactions and interactions, which is pivotal for understanding the public's sentiment on widely viewed content from Police control and a Trump discour.

Link : https://www.youtube.com/watch?v=kuhhT_cBtFU&t=2s
<https://www.youtube.com/watch?v=v3abZ4aAGUU>

Reddit: I harnessed the Reddit API to extract comments from specific threads that sparked debate and controversy. Reddit is a forum known for its community engagement and candid discussions, which can sometimes escalate into toxicity. By focusing on threads with divisive content, I aimed to capture a diverse range of opinions and sentiments, which is crucial for analyzing discourse patterns in an environment that fosters open discussion.

Link :
https://www.reddit.com/r/funny/comments/17r7lh2/was_he_impatient_or_does_he_have_a_point/?onetap_auto=true

Kaggle: The dataset from the "[Jigsaw Toxic Comment Classification Challenge](#)" offered a pre-labeled and structured array of user comments. I selected this particular dataset because it provides a substantial volume of labeled data, which is invaluable for training and benchmarking my machine learning models. The labels cover a range of toxic behaviors, which are essential for a nuanced analysis of online interactions.

BigQuery: Here, I tapped into expansive datasets to process and analyze comments from various online forums, leveraging the computational power of BigQuery to manage the data's scale. Big data systems like BigQuery allow for the processing of massive datasets that would be otherwise challenging to handle. By utilizing this resource, I could include a more comprehensive set of data in my analysis, ensuring that the models I develop are well-informed by a broad spectrum of user interactions.

Link : https://console.cloud.google.com/bigquery?p=bigquery-public-data&d=hacker_news&page=dataset&project=da-bootcamp-2023

Data collection

- **YouTube:** In the "Youtube_WebScraping.ipynb" Jupyter notebook, I've employed a sophisticated technique for web scraping data from YouTube comments. The process begins with setting up the Selenium WebDriver, which is crucial for automated web browsing and interaction with web pages. This setup enables me to programmatically control a web browser, simulating user actions like clicking and scrolling.
Shape : 2238 rows × 4 columns
- **Reddit:** I install PRAW, the Python Reddit API Wrapper, which simplifies the process of accessing Reddit's API. This installation also includes dependencies like update-checker and prawcore. It begins by prompting for Reddit API credentials (client_id and client_secret). These credentials are essential for accessing Reddit's API and are securely entered using Python's getpass module, which hides the client_secret input for security.
Shape : 2219 rows × 4 columns
- **Kaggle:** Like as mentioned on the data and data source section, it was as pretty easy, but one of the most important data for my project. But to get that, it was as just a registration and a simple download.
Shape : 159571 rows × 9 columns
- **BigQuery:** I made sure to have all the necessary libraries by installing google-cloud-bigquery and pandas. I imported necessary modules like os from Python's standard library and bigquery from Google Cloud. My focus here was to write and execute SQL queries to extract data from BigQuery.
Shape : 3000 rows × 4 columns

GDPR

In conducting this analysis, I am acutely aware of the importance of respecting user privacy and adhering to GDPR regulations. To ensure compliance, I've taken several steps:

1. **Data Anonymization:** Any identifiable information extracted during data collection has been anonymized. Usernames and other potential identifiers have been removed or obscured to prevent the possibility of re-identification.
2. **Data Minimization:** I've only collected the data necessary for the analysis. Superfluous details that do not contribute to the project's objectives have been excluded from the datasets.
3. **Consent and Legality:** Where applicable, I've made sure that the data collected is either publicly available or consent for its use in analysis has been obtained, aligning with the legal bases prescribed by the GDPR.
4. **Security Measures:** All datasets are stored securely, with access limited to authorized personnel only, thereby safeguarding the data against unauthorized access or breaches.

By implementing these practices, I ensure that the project not only aligns with legal frameworks but also upholds the principles of ethical data usage.

Data Cleaning and Exploratory Data Analysis

Data Cleaning:

Most of the time there was no need to clean the data in the sense that there were 0 nulls and 0 duplicates, I instead renamed certain columns to harmonize the dataframes with each other. I also removed all special characters from the comments. I added a "comment_length" column to all dataframes. And change the time format of the data columns.

YouTube after : 2212 rows × 5 columns

Reddit after : 2219 rows × 5 columns

Kaggle after : 159571 rows × 9 columns

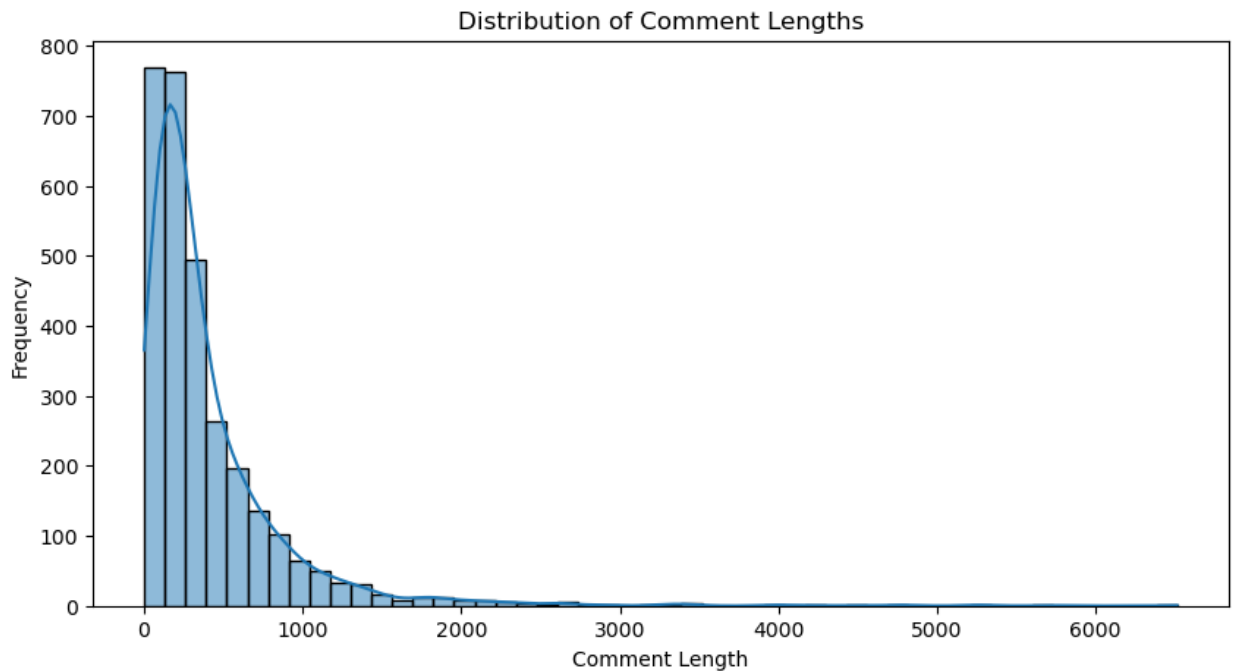
Big Query after : 3000 rows × 5 columns

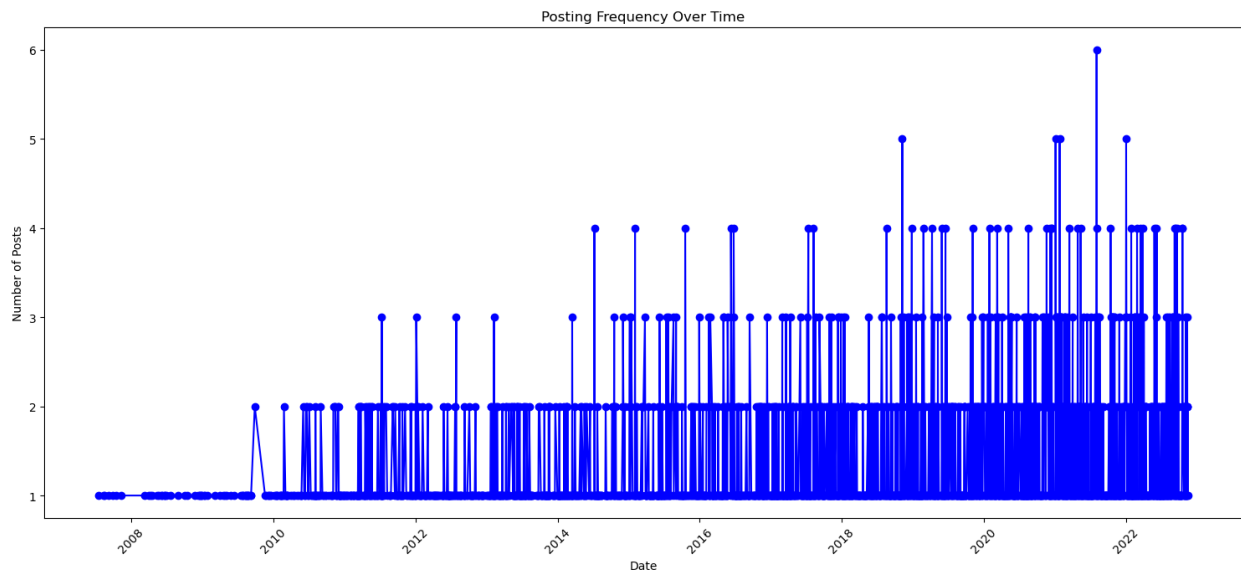
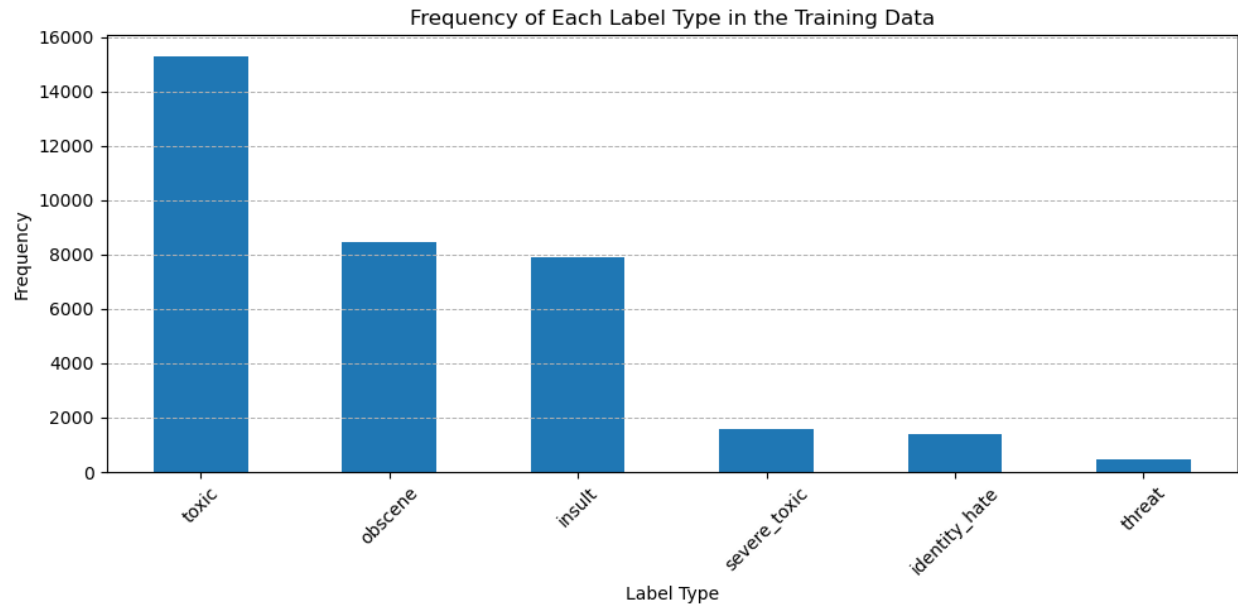
EDA:

1. **Descriptive Statistics:** I calculated descriptive statistics for the 'like' and 'comment_length' columns. It revealed a wide variation in Likes, with a significant standard deviation. Generated a summary: Identified the number of unique authors. Analyzed the time range covered by the dataset, spanning multiple years. Investigated posting frequency to discern patterns over time and analyzed the length of text entries, providing insights into the verbosity and engagement in posts.
2. **Like Distribution:** Using a histogram, I visualized the distribution of comment likes, which displayed a right-skewed pattern indicating that most comments had likes close to zero.
3. **Activity Over Time:** A line graph depicting the number of comments over time showed a decreasing trend in comment activity.
4. **Active Users Analysis:** I identified the most active users by comment count and average like, revealing which users were most engaged in the dataset.
5. **Comment Length Correlation:** By analyzing comment length in relation to likes, I found no strong correlation between the length of a comment and its like.

6. **Word Frequency Analysis:** I generated word clouds for specific users, providing a visual representation of the most frequently used words in their comments.
7. **Special Character Removal:** For UTF-8 compatibility with MySQL, I removed all special characters from the 'comment' column, retaining only alphanumeric characters and spaces. This step was crucial to prevent encoding issues during database import.

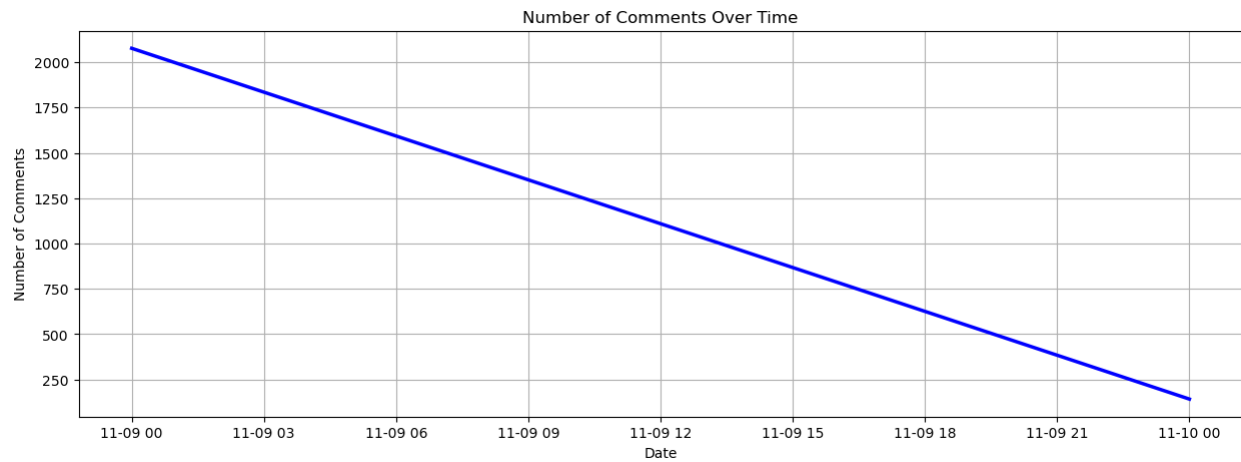
Here's some visualisation examples:





Distribution of toxicity depending of the comment length (%)

	toxic	severe_toxic	obscene	threat	insult	identity_hate
length_category						
(0, 50]	16.854900	2.324639	10.863510	0.425424	9.577108	1.610534
(50, 100]	14.062247	1.629113	7.869995	0.518723	7.533636	1.406225
(100, 200]	9.738177	0.880055	5.121316	0.356408	4.970528	0.841672
(200, 500]	6.921928	0.452371	3.471786	0.203459	3.313781	0.577910
(500, 1000]	5.265800	0.321269	2.610311	0.110436	2.389438	0.496963
(1000, 2000]	4.198430	0.364091	1.934236	0.068267	1.627034	0.352714
(2000, 5000]	8.714134	2.869288	5.313496	0.371945	4.357067	0.956429



[illegible]

Word Cloud for User Global_Lock_2049



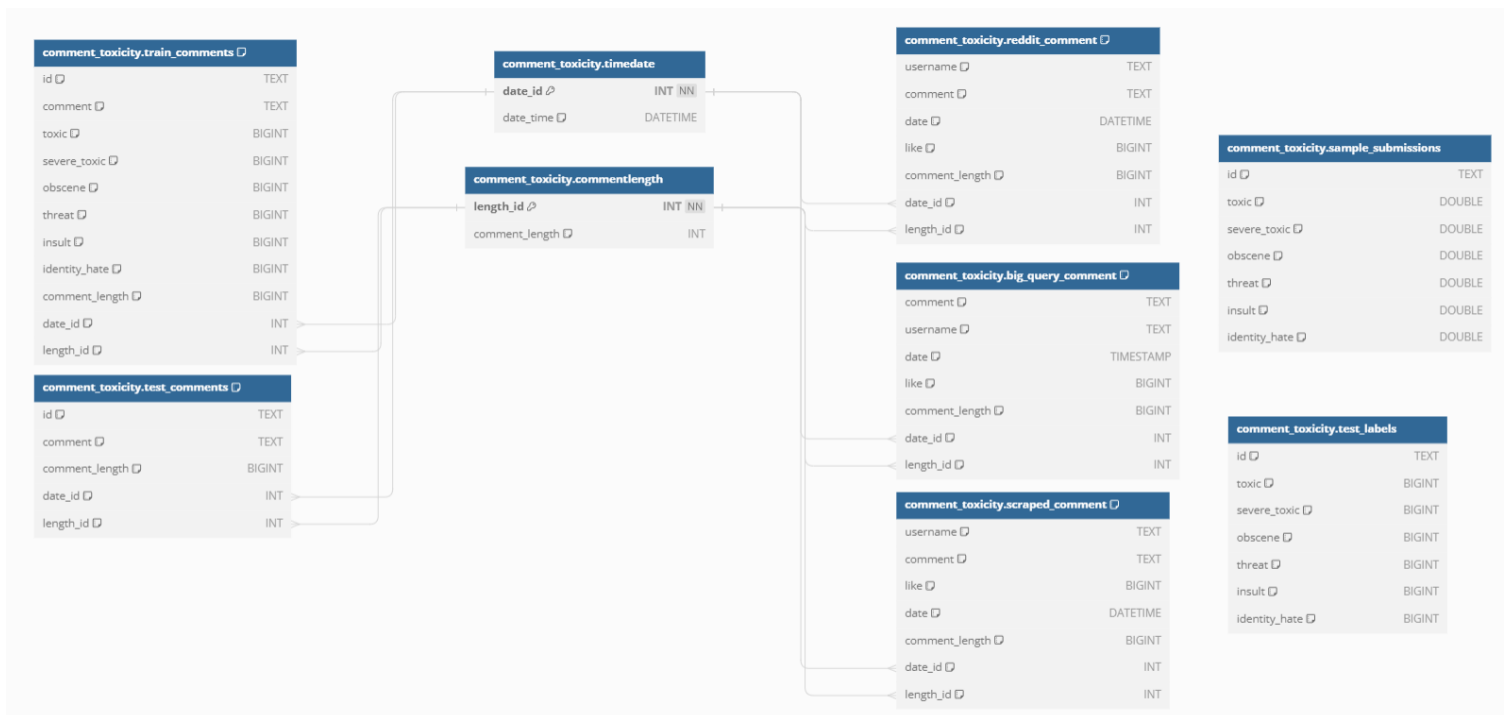
Data base type selection

For this project, MySQL was the database of choice due to its excellent capability to handle structured data, which is abundant in social media commentary. Its relational nature allows me to define clear data schemas, making the organization and retrieval of data more straightforward. Additionally, MySQL's scalability means that as my data grows, the database can efficiently grow with it, which is crucial given the vast amount of comments collected.

MySQL also guarantees the integrity of my data with its transactional support. This is particularly important when integrating multiple data sources, as it ensures consistency throughout the datasets. The security features are top-notch, a non-negotiable aspect when dealing with potentially sensitive user data.

Lastly, the widespread community support and integration with numerous tools and languages streamline the process from data storage to analysis. MySQL's robustness, coupled with its ability to integrate seamlessly into the analytics pipeline, makes it an indispensable tool for my analysis.

Entities. ERD



Database's Queries

```
SELECT td.date_time AS Date, COUNT(*) AS NumberOfComments
FROM TimeDate td
JOIN Scraped_Comment sc ON td.date_id = sc.date_id
GROUP BY td.date_time
ORDER BY td.date_time;
```

```
SELECT AVG(cl.comment_length) AS AverageCommentLength
FROM CommentLength cl
JOIN Scraped_Comment sc ON cl.length_id = sc.length_id;
```

```
SELECT sc.username, COUNT(*) AS NumberOfComments
FROM Scraped_Comment sc
GROUP BY sc.username
ORDER BY NumberOfComments DESC
LIMIT 10;
```

```
SELECT tc.comment, tc.toxic, tc.severe_toxic, tc.obscene, tc.threat, tc.insult, tc.identity_hate
FROM Train_comments tc
WHERE tc.toxic = 1
      AND tc.severe_toxic = 1
      AND tc.obscene = 1
      AND tc.threat = 1
      AND tc.insult = 1
      AND tc.identity_hate = 1
ORDER BY tc.comment
LIMIT 10;
```

```
SELECT AVG(comment_length) AS Average_Length_Of_Highly_Toxic_Comments
FROM train_comments
WHERE toxic = 1
      AND severe_toxic = 1
      AND obscene = 1
      AND threat = 1
      AND insult = 1
      AND identity_hate = 1;
```

Database's Queries – Outputs

Date	NumberOfComments
2022-11-13 10:40:59	380
2022-12-18 10:40:59	95
2023-01-17 10:40:59	18
2023-02-16 10:40:59	14
2023-03-18 10:40:59	26
2023-04-17 10:40:59	16
2023-05-17 10:40:59	26
2023-06-16 10:40:59	21
2023-07-16 10:40:59	27
2023-08-15 10:40:59	20
2023-09-14 10:40:59	41
2023-10-14 10:40:59	47
2023-11-01 10:40:59	4
2023-11-02 10:40:59	1
2023-11-03 10:40:59	1
2023-11-04 10:40:59	2
2023-11-05 10:40:59	1
2023-11-06 10:40:59	5
2023-11-07 10:40:59	2
2023-11-08 10:40:59	8
2023-11-09 10:40:59	3
2023-11-11 10:40:59	1
2023-11-12 14:40:59	1
2023-11-12 20:40:59	1

AverageCommentLength
128.7211

username	NumberOfComments
arifcoco	4
joshuasJR	4
exactteamzis6525	4
gutz1981	3
zmaxpro5681	3
daeneydirusso4069	3
bri9146	3
sfl6307	3
carmenkoening7728	3
stephanidorene	2

comment	toxic	severe_toxic	obscene	threat	insult	identity_hate
I hope your retarded kids get anal rape...	1	1	1	1	1	1
AM GOING TO RAPE YOU IN THE ASS Y...	1	1	1	1	1	1
and your little faggot boy Propol pray pr...	1	1	1	1	1	1
ANYONE WHO SUPPORTS THIS IS FUCK...	1	1	1	1	1	1
Bitch You are a little bitch I fuckin spen...	1	1	1	1	1	1
Eat shit you fucking arse rapping jew fu...	1	1	1	1	1	1
faggot You lil piece of shit I havent van...	1	1	1	1	1	1
FAGGOTS YO FUCKER IT WAS FUCKIN...	1	1	1	1	1	1
Fuck All Asyriac Nation Qamishli belong ...	1	1	1	1	1	1
fuck you honkey why you hatin on black...	1	1	1	1	1	1

Average_Length_Of_Highly_Toxic_Comments
462.5806

Exposing Data via API

```
1 from flask import Flask, request, jsonify, abort
2 import pymysql
3 import os
4 import math
5
6 app = Flask(__name__)
7
8 # Database configuration
9 db_config = {
10     'host': 'localhost',
11     'user': 'root',
12     'password': os.getenv("MYSQL_password"),
13     'db': 'comment_toxicity',
14     'cursorclass': pymysql.cursors.DictCursor
15 }
16
17 # Database connection function
18 def get_db_connection():
19     return pymysql.connect(**db_config)
20
21 @app.route('/')
22 def homepage():
23     ..
24
25
26
27
28
29
30
31
32
33 @app.route('/big_query_comments', methods=['GET'])
34 def get_big_query_comments():
35     username = request.args.get('username')
36     min_like = request.args.get('min_like', type=int)
37     max_like = request.args.get('max_like', type=int)
38     min_length = request.args.get('min_length', type=int)
39     max_length = request.args.get('max_length', type=int)
40     start_date = request.args.get('start_date')
41     end_date = request.args.get('end_date')
42
43     query = "SELECT username, comment, date, `like`, comment_length FROM big_query_comment WHERE 1=1 "
44     params = []
45
46     if username:
47         query += "AND username = %s "
48         params.append(username)
49
50     if min_like is not None:
51         query += "AND `like` >= %s "
52         params.append(min_like)
53
54     if max_like is not None:
55         query += "AND `like` <= %s "
56         params.append(max_like)
57
58     if min_length is not None:
59         query += "AND comment_length >= %s "
60         params.append(min_length)
61
62     if max_length is not None:
63         query += "AND comment_length <= %s "
64         params.append(max_length)
65
66     if start_date:
67         query += "AND date >= %s "
68         params.append(start_date)
69
70     if end_date:
71         query += "AND date <= %s "
72         params.append(end_date)
73
74     conn = get_db_connection()
75     cursor = conn.cursor()
76     cursor.execute(query, tuple(params))
77     comments = cursor.fetchall()
78     conn.close()
79
80     return jsonify(comments)
```

API – Homepage



Welcome to the Comment Toxicity API

This API allows you to access and filter comments based on toxicity metrics and other criteria from various sources like Big Query, Reddit, and scraped data.

Available Endpoints:

/big_query_comments: Fetch comments from Big Query. Can filter by username, likes, comment length, and date.
/reddit_comments: Fetch comments from Reddit. Filters similar to Big Query comments.
/scraped_comments: Fetch scraped comments. Filter options available.
/filtered_train_comments: Fetch comments from training data with specific toxicity metrics.

Example Usage:

To fetch Big Query comments with a specific username:

`http://localhost:5000/big_query_comments?username=johndoe`

To fetch Big Query comments with a specific username:

`http://localhost:5000/big_query_comments?username=johndoe`

To fetch Reddit comments with a comment length greater than 100:

`http://localhost:5000/reddit_comments?min_length=100`

To fetch comments from scraped data with likes greater than or equal to 50:

`http://localhost:5000/scraped_comments?min_like=50`

To fetch training comments marked as toxic and severe_toxic:

`http://localhost:5000/filtered_train_comments?toxic=1&severe_toxic=1`

To fetch Big Query comments within a specific date range:

`http://localhost:5000/big_query_comments?start_date=2021-01-01&end_date=2021-12-31`

Filtering Tips:

Combine multiple query parameters to refine your search.
Use the date format YYYY-MM-DD for date-related queries.
Queries are case-sensitive, especially for usernames.

References

Links:

https://www.youtube.com/watch?v=kuhhT_cBtFU&t=2s

<https://www.youtube.com/watch?v=v3abZ4aAGUU>

[https://www.reddit.com/r/funny/comments/17r7lh2/was he impatient or does he have a point/?onetap auto=true](https://www.reddit.com/r/funny/comments/17r7lh2/was_he_impatient_or_does_he_have_a_point/?onetap_auto=true)

[Jigsaw Toxic Comment Classification Challenge](#)

https://console.cloud.google.com/bigquery?p=bigquery-public-data&d=hacker_news&page=dataset&project=da-bootcamp-2023

TRELLO:

<https://trello.com/b/eUwdijbo/final-project>

GITHUB:

[https://github.com/GuillaumeCapelli/Final Project](https://github.com/GuillaumeCapelli/Final_Project)

THANK YOU