

Hate Speech Filtering Web Extension

- Yulong Cui 190198653
- Arkaitz Zubiaga

Background

- 'Racism' is one of the many expressions of our evolved capacity to live and work in groups.
- Prejudice based on race, skin colour, gender and sexual orientation.
- Most common form of racism is textual based.
- Euro 2020 Final - three England players targeted on social media.
 - More than 2000 racist tweets/posts deleted
 - More than 200+ social media accounts investigated
- Poor mental and physical health related to racism.

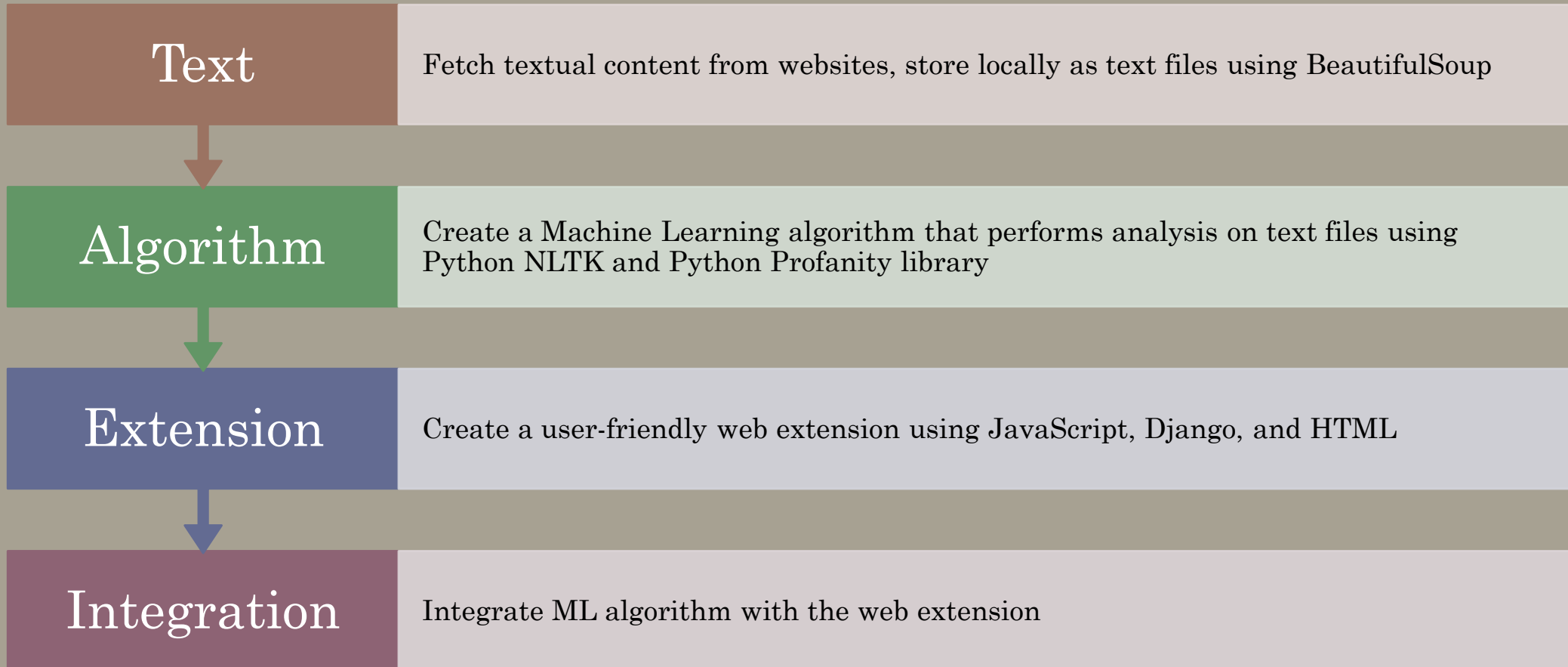
Project Aims

Minimise user's
exposure to online
racism

Reduce reliance on
big social media
companies to filter
out harmful content

User protecting
themselves from
hate speech, be
proactive rather
than reactive

Project Objectives



Computing/Engineering problem

- Filtering textual based content requires:
 - Natural Language Processing (NLP)
 - Variations of words
 - Combinations of phrases
 - Understanding in context
 - Machine Learning algorithm
 - Pattern recognizing
 - ‘Training data’ oriented

Literature Review Findings



Python & related
libraries



Machine Learning



Existing systems

Python & related libraries

Python programming language

- Object-oriented
- Dynamic typing, binding
- Extensive selection of Machine Learning libraries

Natural Language Processing (NLTK libraries)

- Text processing libraries:
 - Classification
 - Tokenization
 - Stemming
 - Tagging
 - Parsing
 - Semantic reasoning

Profanity filter

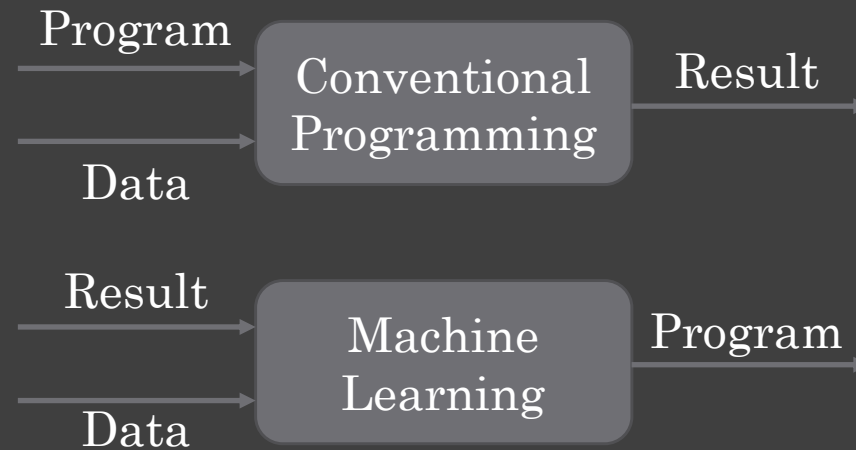
- Word censoring
- Deep analysis
- RESTful web service

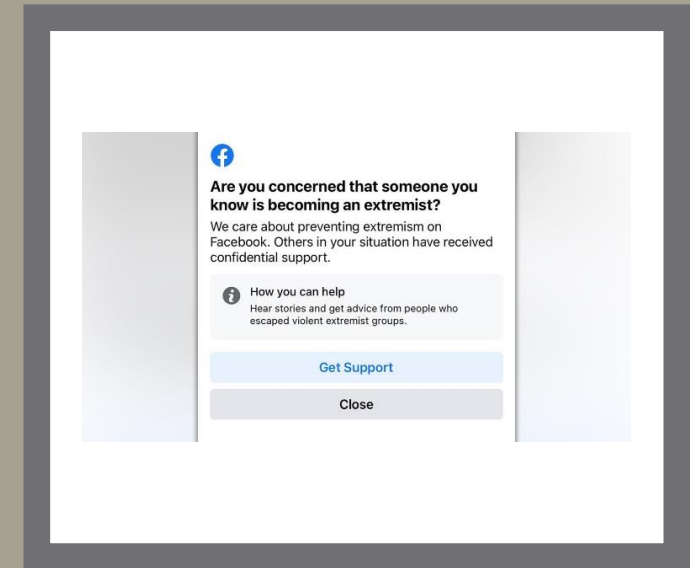
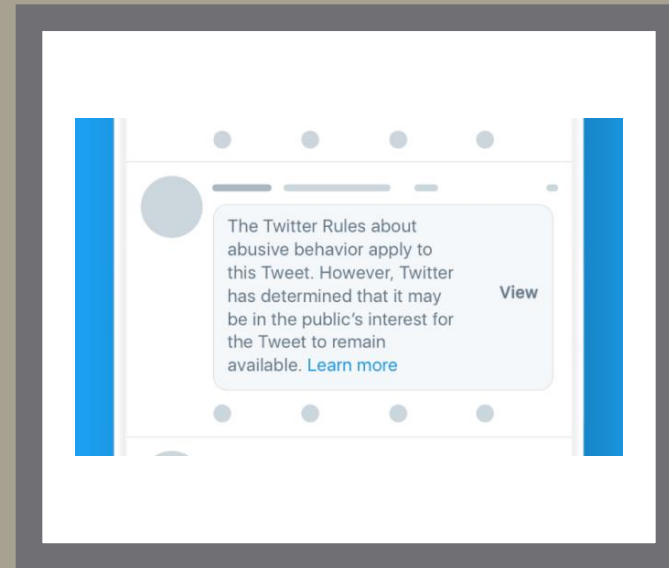
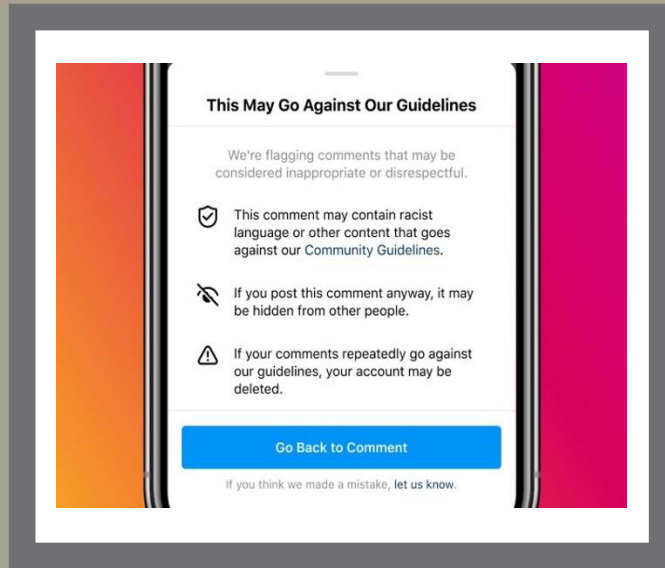
Beautiful Soup

- Utilizes HTML/XML parser
- Web scraper

Machine Learning


- ML algorithms build a mathematical model/program using pre-written algorithm depending on 'training data'.
- Conventional programming outputs results depending on user written logic/programs.






Existing systems

Meta (Facebook) social media platforms' content filtering system


Hate Block
 Offered by: Hate Block
 ★★★★★ 0 | [Social & Communication](#) | 5 users

[Add to Chrome](#)

[Overview](#)
[Privacy practices](#)
[Reviews](#)
[Related](#)



Overview

Compatible with your device

Hate speech has been at an all-time high on the internet. This amounts to not just harassment, but also creates misinformation & anger.

Hate speech has been at an all-time high on the internet. And this amounts to not just harassment and mental disturbance, but also creates misinformation and anger. As India will see a rise of this hate speech during the 2019 Elections, it is time to give the citizens a way to #TurnOffTheHate.

Additional Information

[Website](#)
[Report abuse](#)


Version
 0.0.10

Updated
 27 April 2019

Size
 71.9KiB

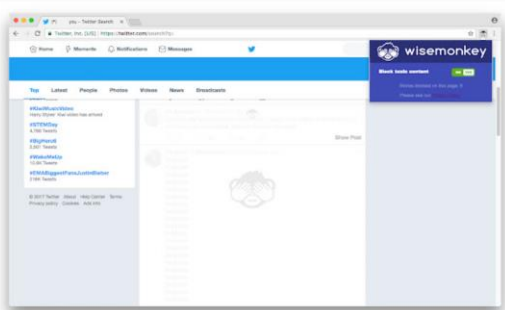
Language
 English

Developer
[Contact the developer](#)


WiseMonkey
 Offered by: wisemonkey.io
 ★★★★★ 4 | [Social & Communication](#) | 18 users

[Add to Chrome](#)

[Overview](#)
[Privacy practices](#)
[Reviews](#)
[Related](#)



Overview

Compatible with your device

Automatically block cyberbullying, harassment, hate speech on Twitter and Facebook using machine learning and data science

WiseMonkey uses new artificial intelligence and NLP technology to filter cyberbullying, harassment, and hate speech from your Twitter and Facebook news feeds. Learn more at <http://wisemonkey.io>.

Additional Information

[Website](#)
[Report abuse](#)

Version
 0.17.4

Updated
 8 November 2017

Size
 50.1KiB

Language
 English

Developer
[Contact the developer](#)

Existing systems

Similar Chrome web browser extensions

GUI Design

-Extension Interface



Hate Speech Filtering

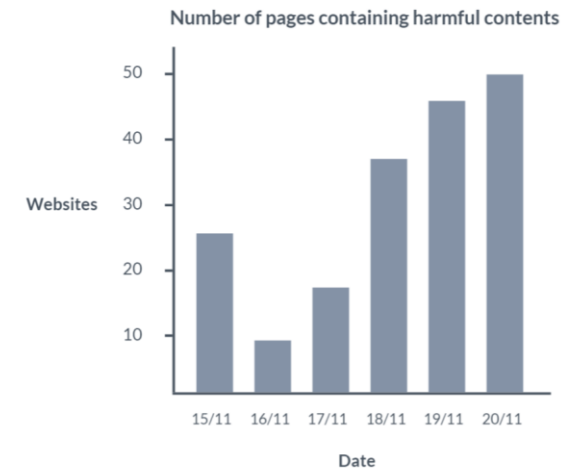
Filter texts on

This website:

<https://twitter.com/sampleUserName/status/123456>

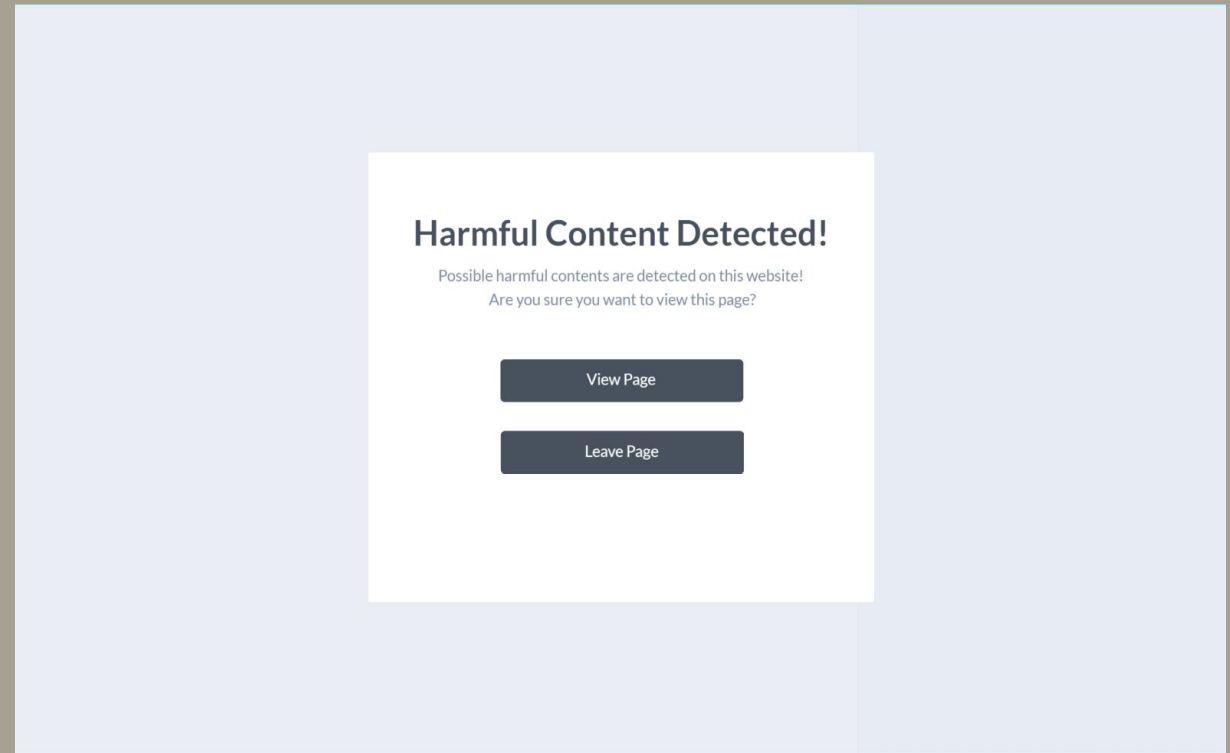


Number of websites contains harmful contents so far:



GUI Design

-Warning splash screen



Evaluation

- Text scraper
 - Social media websites
 - Tweets/posts
 - Replies/comments
 - Published documents
- Machine Learning algorithm
 - Training data
 - Existing results
- Web extension
 - Developer testing
 - Sample user testing

Project Planning

[illegible]

Risk Register

Risk Description	Likelihood rating	Impact rating	Impact	Preventative actions
Poor time management	Medium	High	Affect entire project timeframe, resulting in unfinished final product or a product with less standards	Utilise Gantt chart to track progress regularly, ensure sufficient progress is made on time or ahead of the time
Personal illness and unforeseen events	Low	Medium	Cause delays in final product	Ensure healthy personal schedule and diet, protect mental health
Lack of technical expertise	Medium	High	Cause delays in final product and lower standards in final product	Ensure regular contact with advisor, seek help from experienced peers if necessary
Loss of code/progress	Low	High	Loss of large portions of work if forgot to save progress/computer problems	Use cloud services such as GitHub and OneDrive for code backups
Project development becomes too large to finish before deadline	Low	Medium	Project overextends to unnecessary territories causing loss focus on main functionalities	Stick with original plan for project, minor changes can be done by myself, consult experienced peers and advisor for any major changes to project

References

1. Brooks, R., 2021. The Origins of Racism. [online] The Conversation. Available at: <<https://theconversation.com/the-origins-of-racism-8321>> [Accessed 26 November 2021].
2. Perrigo, B., 2021. *Twitter Offers Transparency on Racist Abuse, But Few Solutions*. [online] Time. Available at: <<https://time.com/6089289/twitter-racist-abuse-anonymity/>> [Accessed 26 November 2021].
3. Lewsley, J., (2020). The effects of racism on health and mental health [online]. *Medical and health information*. Available at: <https://www.medicalnewstoday.com/articles/effects-of-racism#adults> [Accessed 26 November 2021].
4. What is python? Executive summary [online]. *Python.org*. [Viewed 27 November 2021]. Available from: <https://www.python.org/doc/essays/blurb/>
5. Python AI: why is python so good for machine learning? [online]. *Custom Software Development | Netguru*. [Viewed 27 November 2021]. Available from: <https://www.netguru.com/blog/python-machine-learning>
6. NLTK : natural language toolkit [online]. (no date). *NLTK : Natural Language Toolkit*. [Viewed 27 November 2021]. Available from: <https://www.nltk.org/>
7. Profanity-filter [online]. (no date). *PyPI*. [Viewed 27 November 2021]. Available from: <https://pypi.org/project/profanity-filter/#deep-analysis>
8. BeautifulSoup4 [online]. *PyPI*. [Viewed 27 November 2021]. Available from: <https://pypi.org/project/beautifulsoup4/>
9. How different are conventional programming and machine learning? - [online]. *Tech Blogger*. [Viewed 28 November 2021]. Available from: <https://contenteratechspace.com/how-different-are-conventional-programming-and-machine-learning/>