# tbs Think & Create

Toulouse
Business School

## Data Project

**Build Your Own Hate speech detection Engine**

January 2021

**www.tbs-education.fr**

# Definition

- **Encyclopedia of the American Constitution:** "Hate speech is speech that attacks a person or group on the basis of attributes such as race, religion, ethnic origin, national origin, sex, disability, sexual orientation, or gender identity."

- **Facebook:** "We define hate speech as a direct attack on people based on what we call protected characteristics—race, ethnicity, national origin, religious affiliation, sexual orientation, caste, sex, gender, gender identity, and serious disease or disability.

- **Twitter:** "Hateful conduct: You may not promote violence against or directly attack or threaten other people on the basis of race, ethnicity, national origin, sexual orientation, gender, gender identity, religious affiliation, age, disability, or serious disease."

# Definition

- Hate speech is to incite violence or hate

- Hate speech is to attack or diminish

- Hate speech has specific targets

- Whether humor can be considered hate speech

# Motivation

- Hate crimes are nothing new in society;

- Social media have begun playing a major role in hate crimes;

- Suspects in several recent hate-related terror attacks had an extensive social media history of hate-related posts;

- Social media could contribute to their radicalization;

- Many online forums such as Facebook, YouTube, and Twitter consider hate speech harmful,

# Your mission

- You are part of the TWITTER team of data scientists.

- This company offers  microblogging and social networking service on which users post and interact with messages known as "tweets".

- She would like to develop a new application able to detect hate speech and offensive messages.

# The objective

- **For a given tweet/retweet, your application should be able to classify the comment as a hate message, an offensive message or neither hate speech nor offensive;**

- **"Rshiny" should be used to convert your R codes to an Application;**

- **Git and Github should be used as collaborative working tools,**

# Possible Datasets

- **The dataset1 is stored as a CSV. Each data file contains 5 columns:**

  - **count** = number of experts who coded each tweet (min is 3, sometimes more experts coded a tweet when judgments were determined to be unreliable).

  - **hate_speech** = number of experts who judged the tweet to be hate speech.

  - **offensive_language** = number of experts who judged the tweet to be offensive.

  - **neither** = number of experts who judged the tweet to be neither offensive nor non-offensive.

  - **class** = class label for majority of experts: 0 - hate speech 1 - offensive language 2 - neither

- **The dataset2 contain a non-exhaustive list of Abuse words**

  - **404 abuse terms** used in:

    - "Twits, Twats and Twaddle: Trends in Online Abuse towards UK Politicians", ICWSM 2018,

    - "Online abuse of uk mps in 2015 and 2017: Perpetrators, targets, and topics", extended version on arXiv, 2018.

  - **388 abuse terms** used in:

    - "Online Abuse of UK MPs from 2015 to 2019: Working Paper", arXiv working paper.

# Possible Datasets

- **The dataset3 contain a list of 1528 annotated comments from Fox News website.**

    - **435 of them are labeled as hateful;**

    - **1093 of them are labeled as non-hateful.**

- The dataset4 contain a list of abusive microposts that were considered to be explicitly abusive

  - The first line of each micropost represents the class label which is either "abusive" or "notAbusive"

- **No restriction on using other datasets !!!!!**

# Methodology

- **Text mining ?**

- **Logistic regression ?**

- **Naïve Bayes ?**

- **Decision trees ?**

- **Other Machine Learning Technics ?**

# Methodology

- **Some Related papers:**

  - **Automated Hate Speech Detection and the Problem of Offensive Language, Thomas Davidson et al. (2017);**

  - **Hate Speech Detection: A Solved Problem? The Challenging Case of Long Tail on Twitter, Ziqi Zhang et al. (2018);**

  - **Hate speech detection: Challenges and solutions, Sean MacAvaney et al. (2019)**

# Evaluation criteria: Groups of 05 students

- Evaluation criteria
  - Understanding of business context (identification of key points, understanding of sector) Mastery & Pertinence of statistical analysis (descriptive, modeling);
  - Link between data / analysis / business recommendations (logical flow, pertinence, completeness);
  - Business concepts (thoroughness, pertinence);
  - Thoroughness and quality of business recommendations (pertinence, relevance, professionalism, originality);
  - Understanding of big data (mastery of big data terms, confidence when using big data concepts);
  - Visuals (professionalism, slides support well the main arguments of the presentation, appropriate content);
  - Delivery (clear and logical organization, effective introduction and conclusion, creativity, transition between speakers, oral communication skills, eye contact);
  - Q&A session (ability to answer questions);
  - Report: Quality of data analysis,

# Evaluation criteria: Groups of 04 students

- **Very important:**

  - **Try to use materials you learned during all the Program**

  - **We want to be impressed !!!**

  - **Presentation during the last session:**

    - 20 minutes per groups;

    - Presentation in front of jury members.

# Use Github

- A collaborative tool to work together on projects;

- A code hosting platform for version control and collaboration;

- Basic services are free of charge;

- More advanced professional and enterprise services are commercial;

- Free GitHub accounts are commonly used to host open-source projects;

- Free plan: unlimited collaborators, private repositories restricted to 2,000 minutes of GitHub Actions per month;

- 40 million users and more than 190 million repositories;

- The largest host of source code in the world;

- Since October 26, 2018: Acquired by Microsoft.

# Get started with Github

- Version control

  - The management system that manages the changes that you made

    in the project: addind new files, modifying older files, etc.

  - Every time you make a change on your project a different version

    is created and saved (called a snapshot);

  - All the versions are kept,

- **Version control**



> Version control is the management of changes to documents, computer programs, large web sites, and other collections of information.

> These changes are usually termed as "versions".

# Get started with Github

- **Why a Version control ?**
  - **You are always notified to changes on files;**
  - **Avoid change conflicts;**
  - **Facilitates collaboration**

- **Version control System Tools**

# Get started with Github

- **Version control System Tools**

- **Git and Github**
  - **Git**: management version control tool; help to create local repository, pull data from the central server and push local data to the server;
  - **Github**: central repository; code hosting platform for version control collaboration.

- **Overview of Git**
    - A **repository** is a directory or a storage space where your project can live; can be local on your computer or a storage space on Github; It keep every thing related to your project;

Git is a Distributed Version Control tool that supports distributed non-linear workflows by providing data assurance for developing quality software.

- **Central and local repository**

- **Git Operations and Commands**

# Get started with Github

- **Creating Repositories**

GitHub

Create your Central Repository on GitHub

git

**git init**

Install Git on your local machine and use "git init" to create your local repository.

**git clone**

OR

Download or clone your repository from GitHub.

# Get started with Github

- **Creating a central Repository on Github**

  - **Create an account on** https://github.com/

  - Create a new repository

    - In the upper right corner, next to your avatar or identicon, click and then select

      New repository.

    - Name your repository.

    - Write a short description.

    - Select Initialize this repository with a README.

# Get started with Github

- **Creating a central Repository on Github**

# Get started with Github

- **Creating a central Repository on Github**

# Get started with Github

- **Creating a local Repository on my local machine**

  - **Install Git on your computer**

    - **For window:** https://git-scm.com/download/win

    - **For Mac:** https://git-scm.com/download/mac

  - **Create a project file in your "Local Disk (C)"**

# Get started with Github

- **Open the local created file and right click**

- **Click on "Git Bach Here"**

- **A Git Bach Emulator is opened**

# Get started with Github

- **Open the local created file and right click**
- **Click on "Git Bach Here"**
- **A Git Bach Emulator is opened**
- **Your commands are done in this bach**

# Get started with Github

- **To create your local repository type in this bach: "git init" and press "enter"**

# Get started with Github

- **A ".git" folder has been created with all information, objects, etc**

- **Link the local and the central repositories**



- Use '**git add origin <link>**' to add remote repo.
- Pull files with '**git pull**'
- Push your own changes into central repo with '**git push**'

# Get started with Github

- **Link the local and the central repositories**
  - **Add the remote repository as your origin by typing into the bach:**
    - git remote add origin "http url of your Github "

# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Add new files to your index by typing into the bach the command:**

    **git add filename1 filename2**

# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Commit the changes into the bach the command:**

    **git commit -m "your message"**

# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Generate an SSH public key from your bach:**

    **ssh-keygen**

# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Print your SSH public key from your bach:**

# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Add your SSH public key in your Github account:**
    **Settings>SSH and GPG Keys>New SSH key**
  - **Give a name to your key, paste the key code, and add a SSH Key**
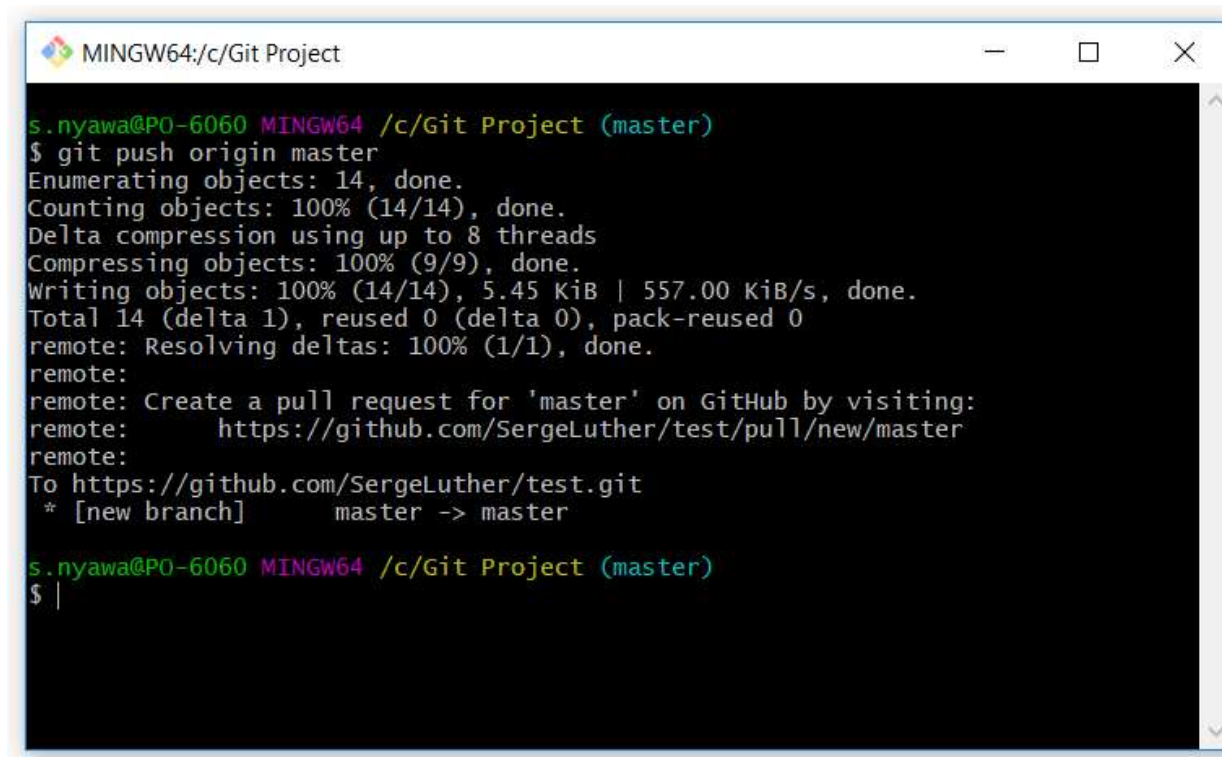
# Get started with Github

- **Pushing files from the local repository to the central repository**
  - **Authentify your key from your bach:**

```
$ ssh -T git@github.com
The authenticity of host 'github.com (140.82.121.3)' can't be established.
RSA key fingerprint is SHA256:nThbg6kXUpJWGl7E1IGOCspRomTxdCARLviKw6E5SY8.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'github.com,140.82.121.3' (RSA) to the list of known
hosts.
Hi SergeLuther! You've successfully authenticated, but GitHub does not provide s
hell access.

s.nyawa@PO-6060 MINGW64 /c/Git Project (master)
$ |
```

- **Pushing files from the local repository to the central repository**
  - **Push the files from your bach:**

# tbs Think & Create

## Toulouse Business School

www.tbs-education.fr