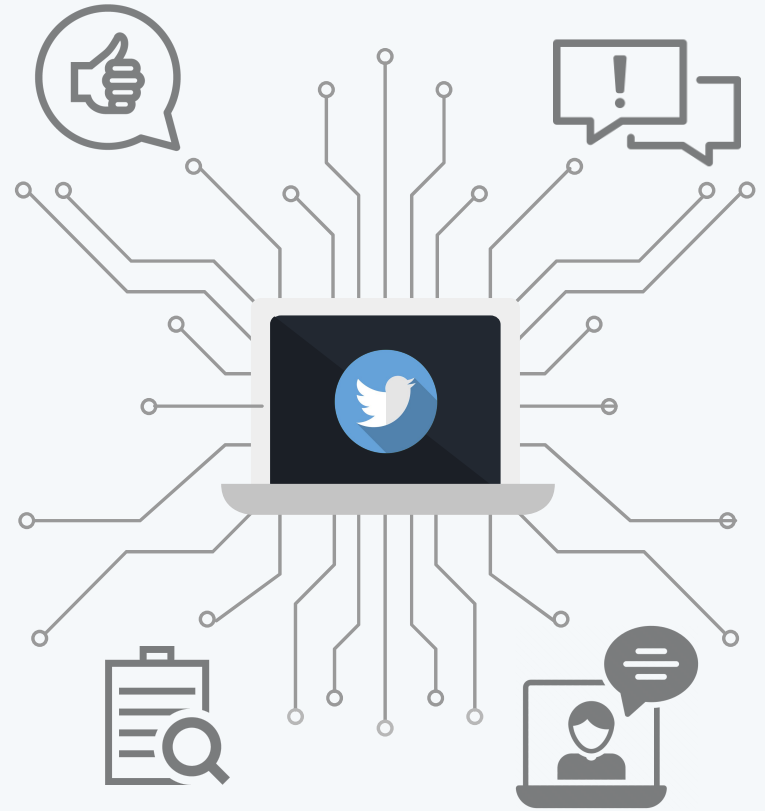# Twitter Hate Speech Detection

*Can Content Moderation be Automated?*
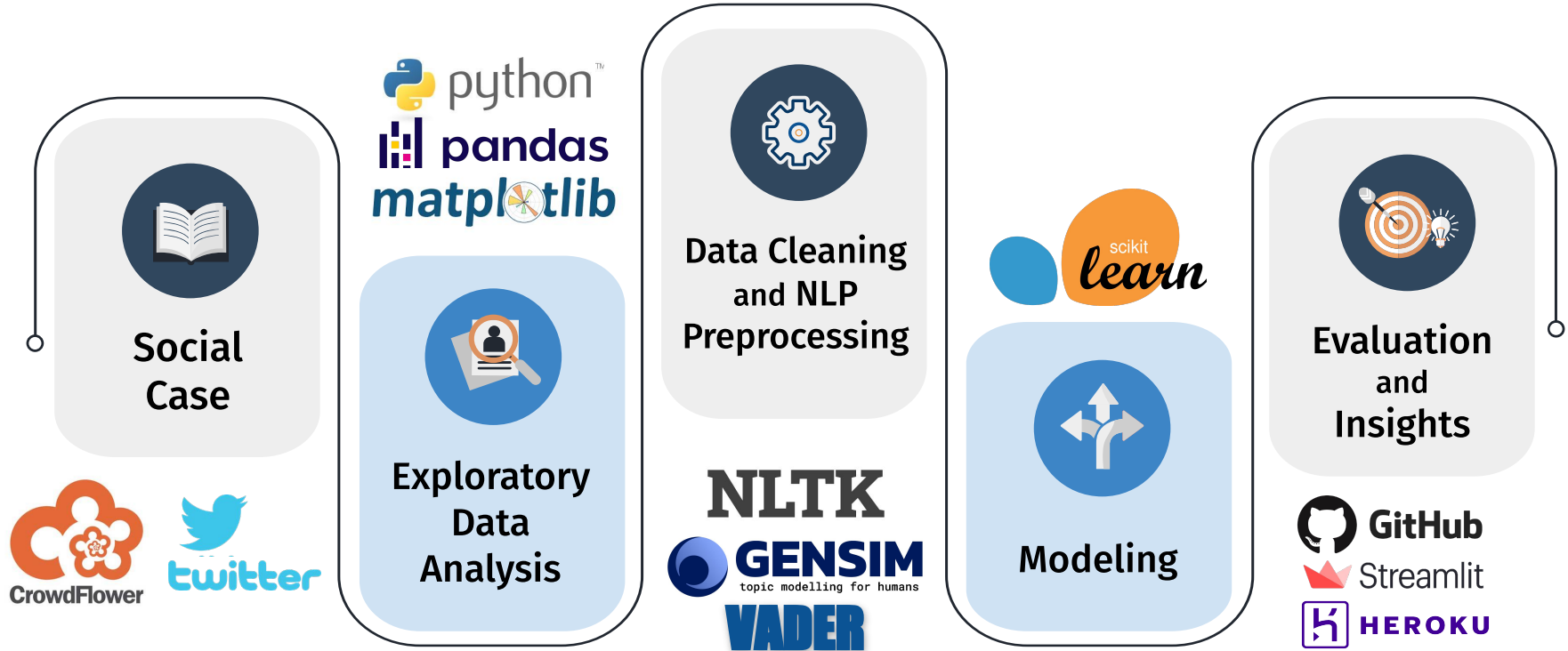
Flatiron School Capstone
**Sidney Kung**

# The Problem of Human Content Moderation
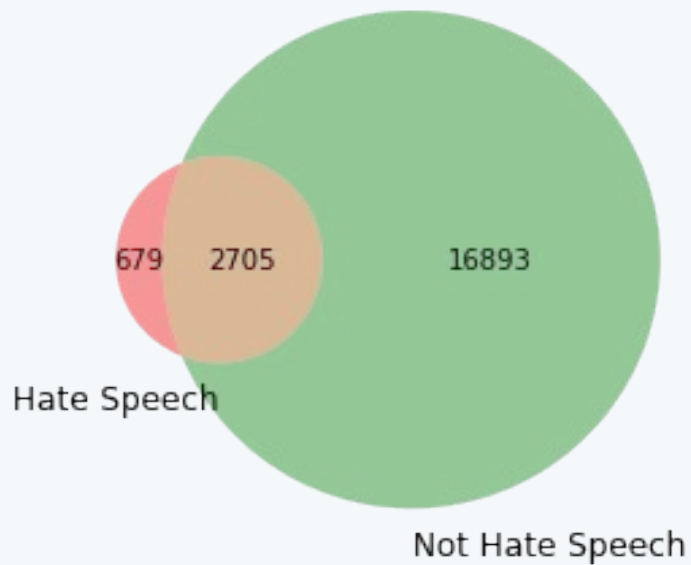
- **Every major tech company** uses third-party contractors

- **Automating** this process could **reduce labor exploitation**

- What is **Hate Speech?**
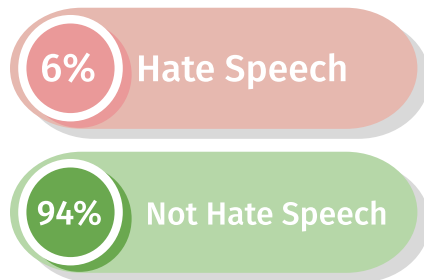
# CRISP-DM Process

**Unique Words per Label**

679    2705      16893

Hate Speech

Not Hate Speech

# Data Understanding

Sourced from 2017 Cornell University **research study**.

**24,802** Tweets

**20,277** Word Vocabulary

**6%** Hate Speech

**94%** Not Hate Speech

# Data Analysis

**1** What are the **linguistic differences** between hate speech and offensive language?
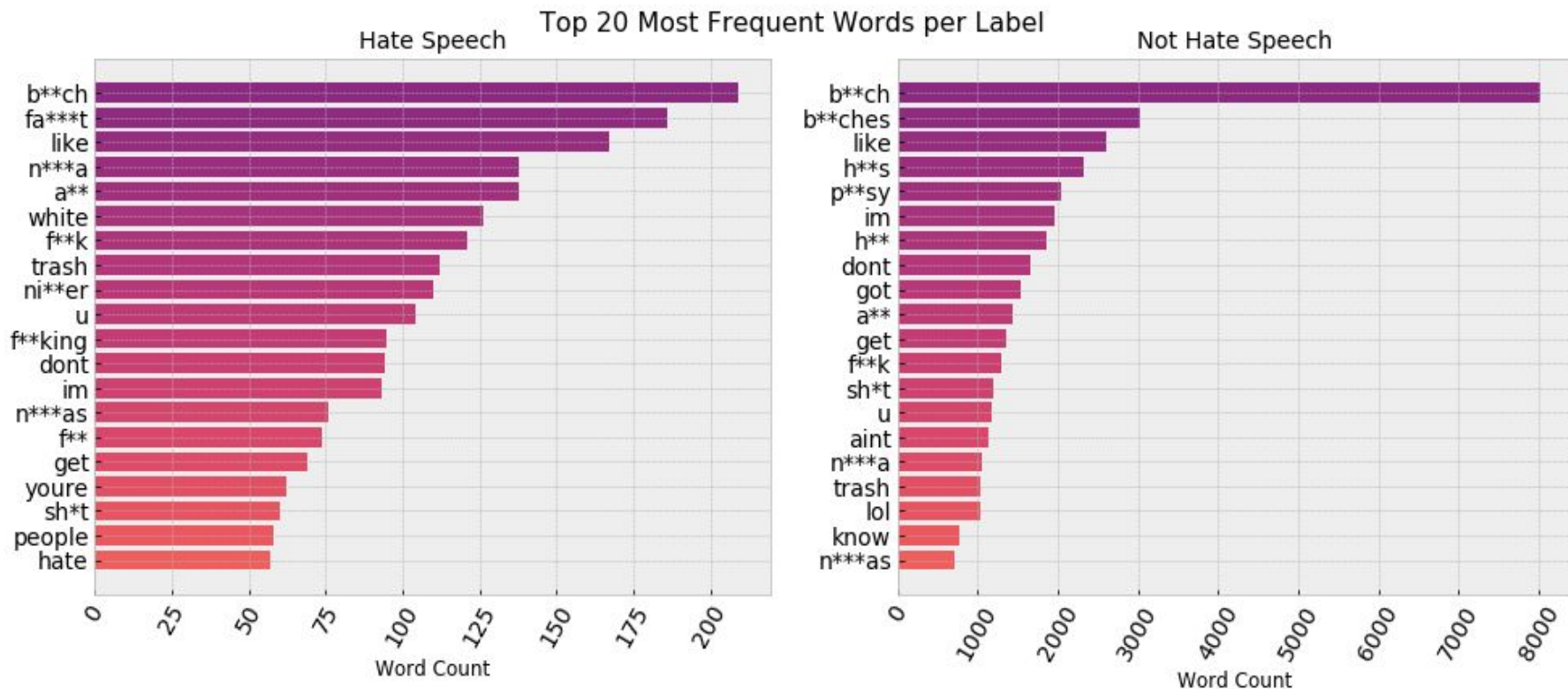
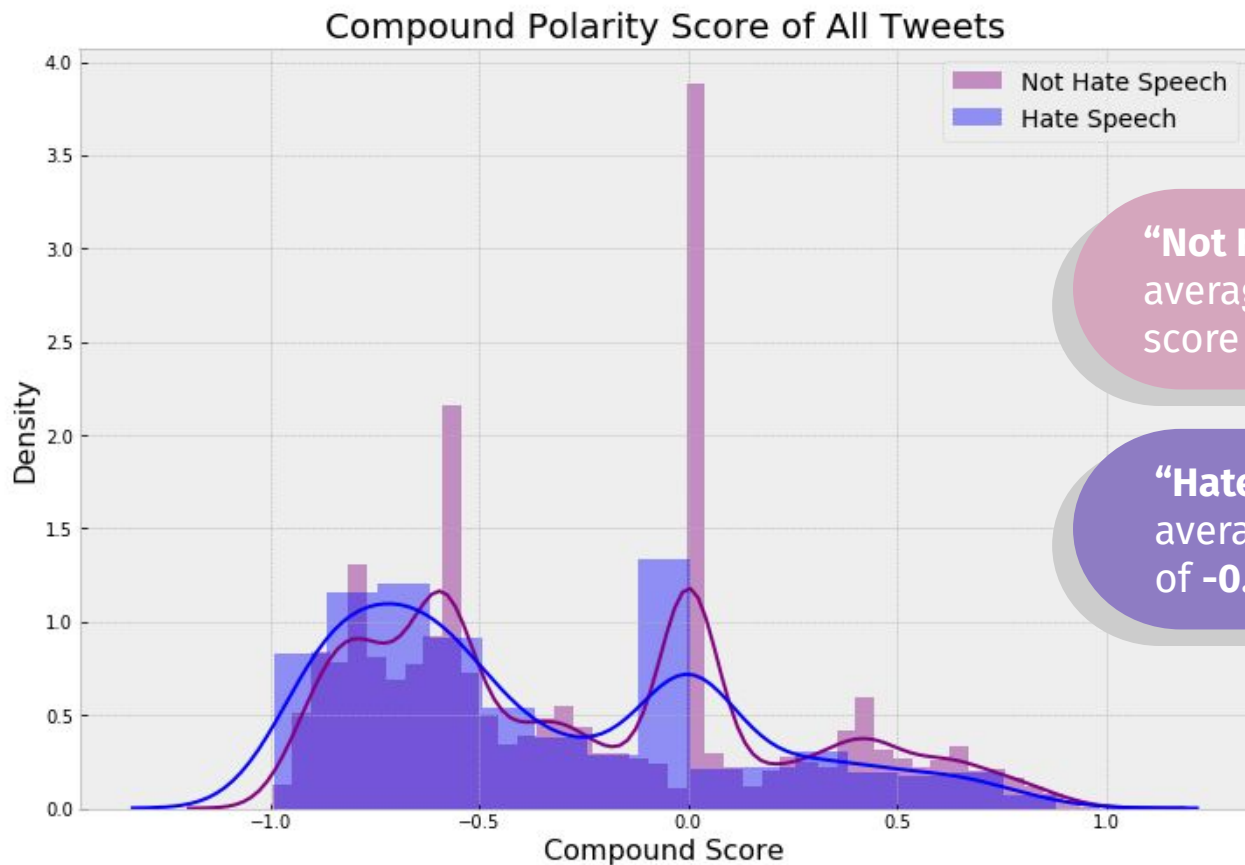**2** What are the **popular hashtags** of each tweet type?

**3** What is the **overall polarity** of the tweets?

# What are the **linguistic differences** between hate speech and offensive language?


Top 20 Most Frequent Words per Label

# What are the **popular hashtags** of each tweet type?

# What is the **overall polarity** of the tweets?



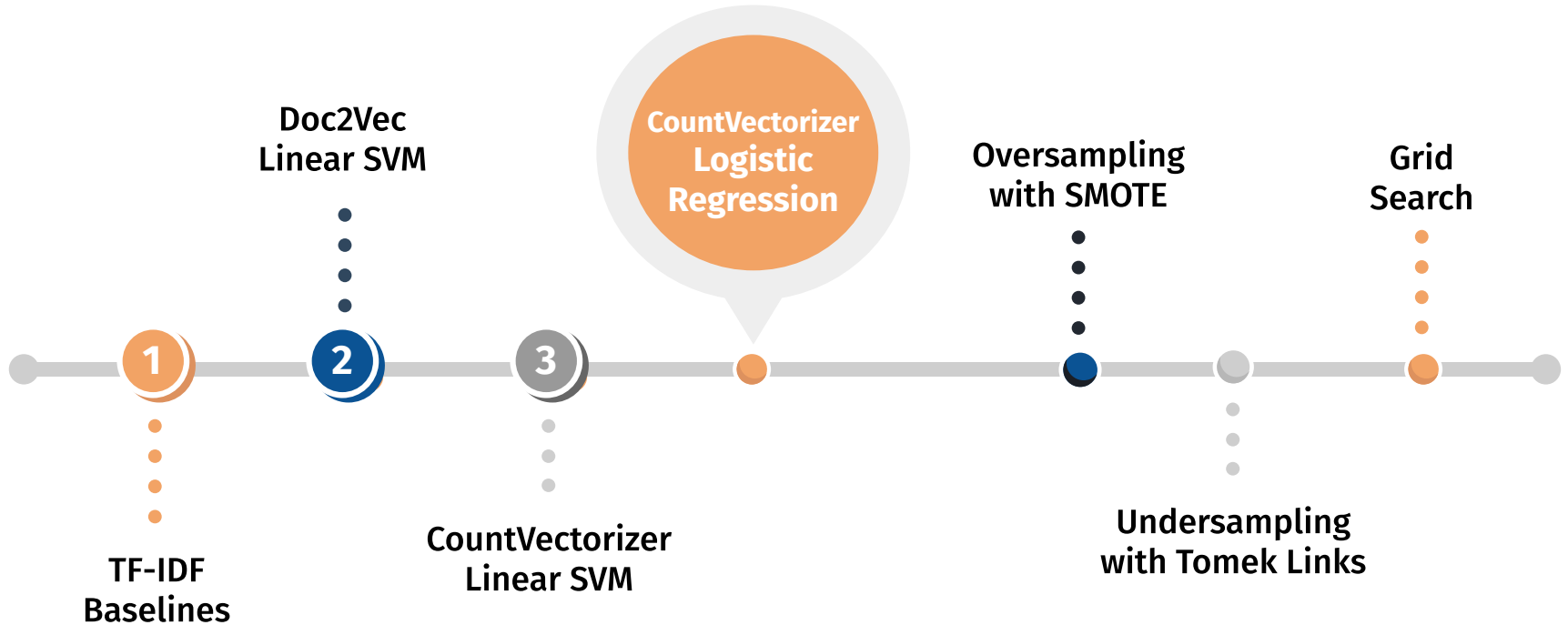Compound Polarity Score of All Tweets

**"Not Hate Speech" tweets:** average compound score of **-0.263**

**"Hate Speech" tweets:** average compound score of **-0.363**

# Modeling Process

## Model Deployment

**hate-speech-predictor. herokuapp.com**

## Is Your Tweet Considered Hate Speech?

*Please note that this prediction is based on how the model was trained, so it may not be an accurate representation.*

Enter Tweet

0/280

**Prediction:**

## For More Information

**Check out the project repository here.**

Contact Sidney Kung via sidneyjkung@gmail.com.

**Let's Connect!**

LinkedIn | Github | Medium | Twitter

# Thank You!

**GitHub Repository**
github.com/sidneykung/
twitter_hate_speech_detection

**Web App on Heroku**
hate-speech-predictor.
herokuapp.com

SidneyKung.com

SidneyJKung@gmail.com

linkedin.com/in/sidneykung

@Sidney_K98