

# Twitter Hate Speech Detection

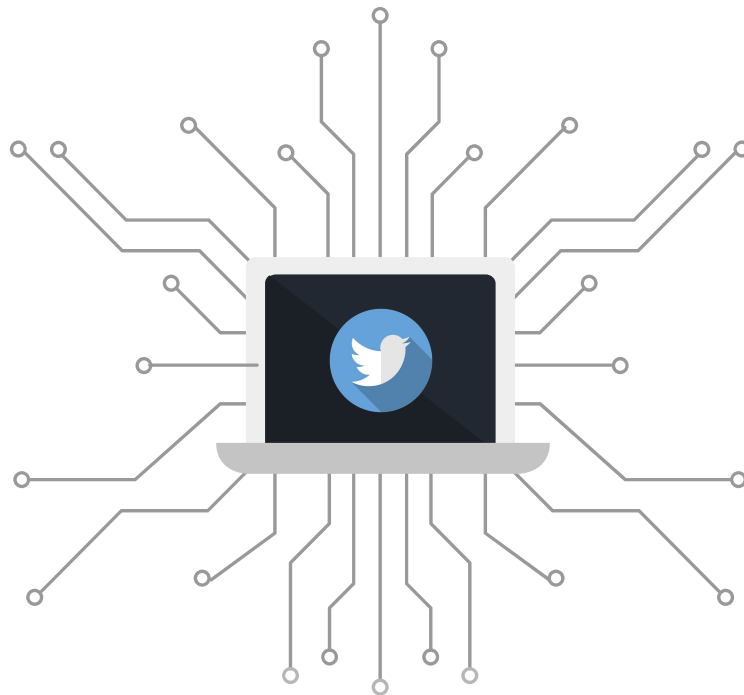
*Can Content Moderation  
be Automated?*

Flatiron School Capstone  
**Sidney Kung**

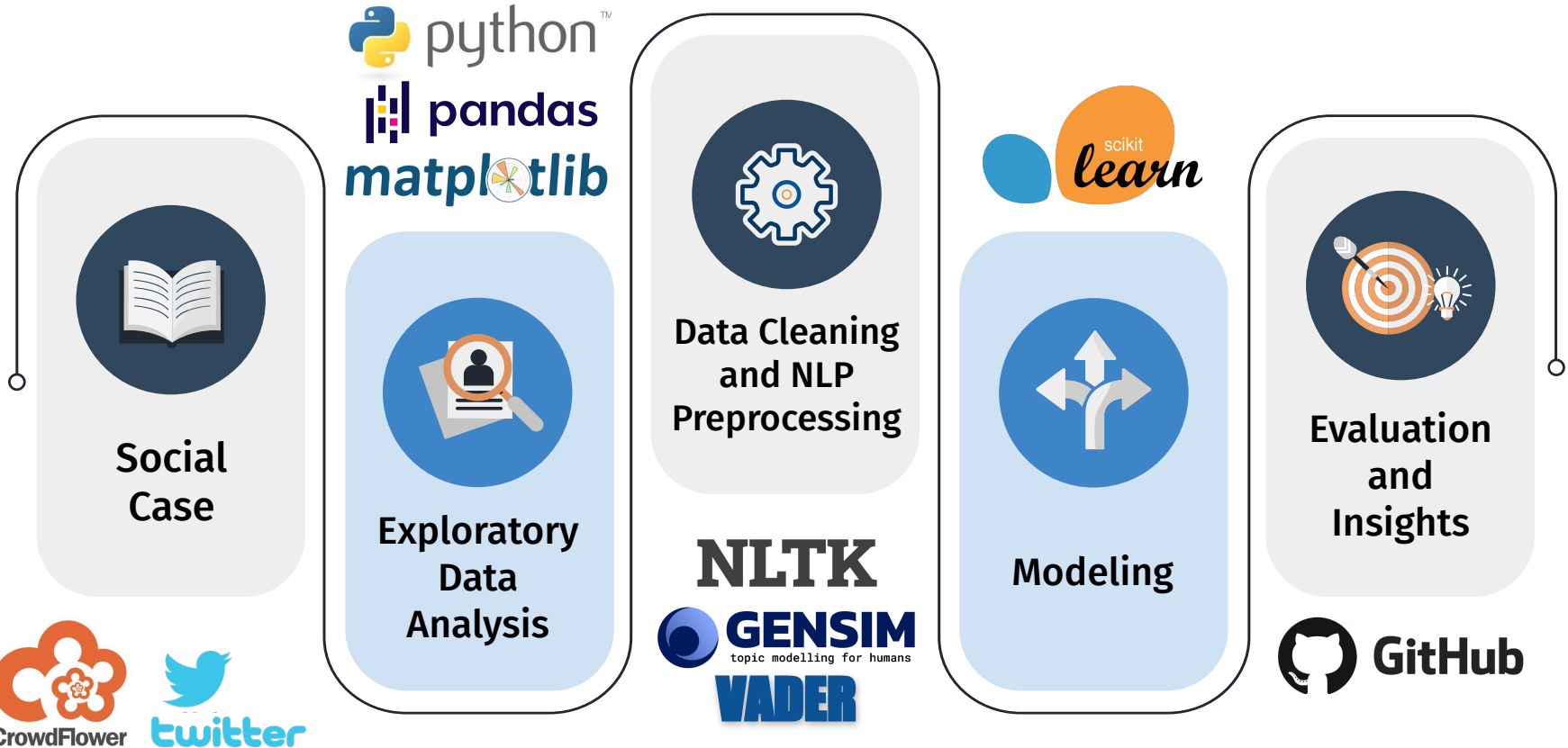


# Overview

- 2019 Verge article exposed Cognizant, a former Facebook **content moderation contractor**
- **Automating hate speech detection** could **reduce labor exploitation** and other human rights violations
- **Hate Speech** is defined as abusive or threatening speech that expresses prejudice against a particular group.



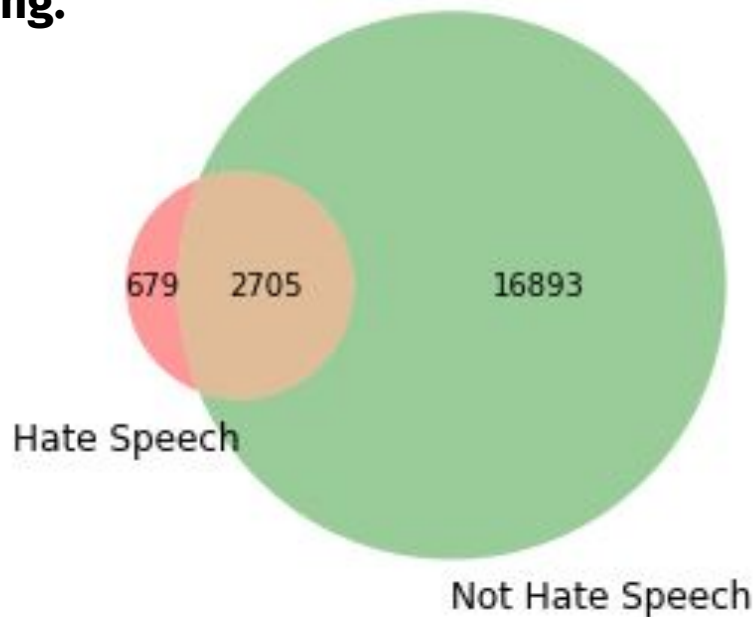
# CRISP-DM Process



# Data Understanding

Sourced from Cornell University **research study**.  
**Labeled by** CrowdFlower **majority-rules voting**.

- **24,802** Tweets
- Vocabulary of **20,277 unique words**
- Binary Classification
  - 6% Hate Speech
  - 94% Not Hate Speech
- Evaluation Metric: **F1 Score**



# Business Questions

1

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



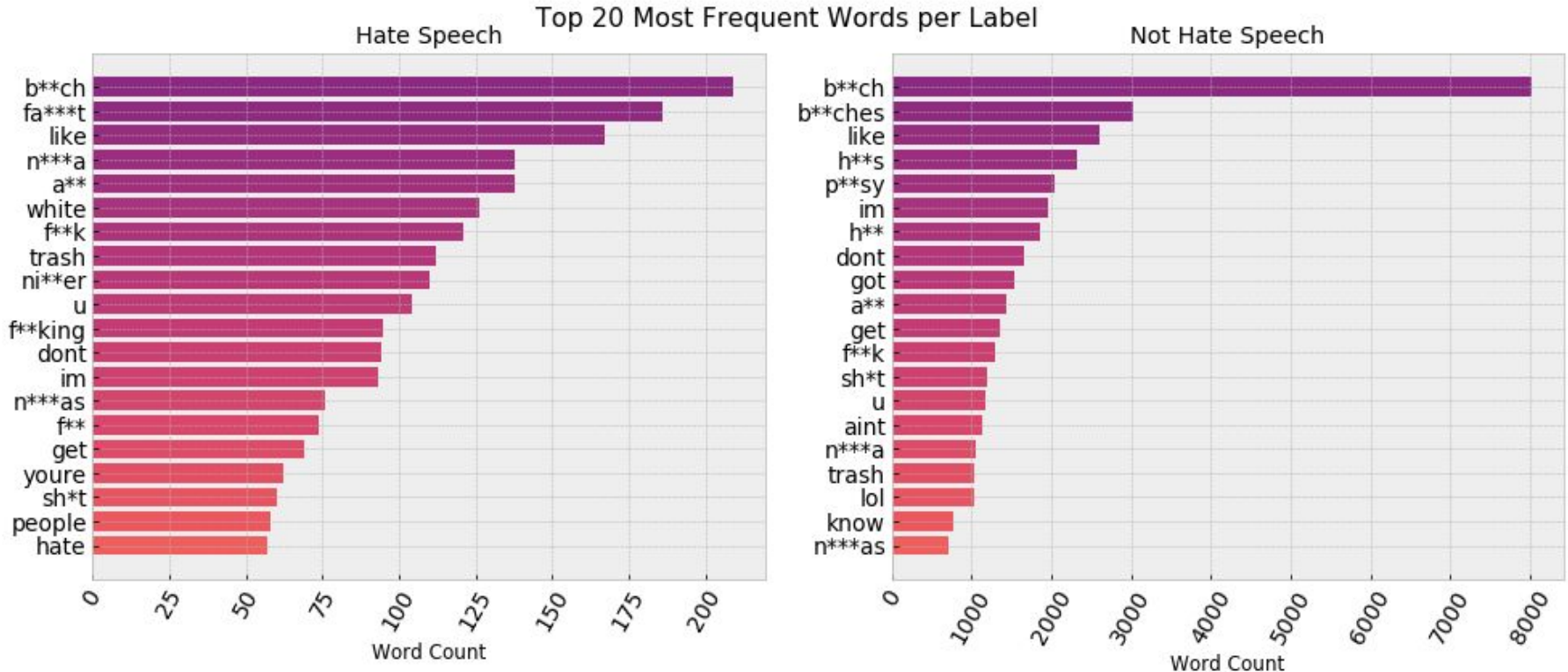
2

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

3

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

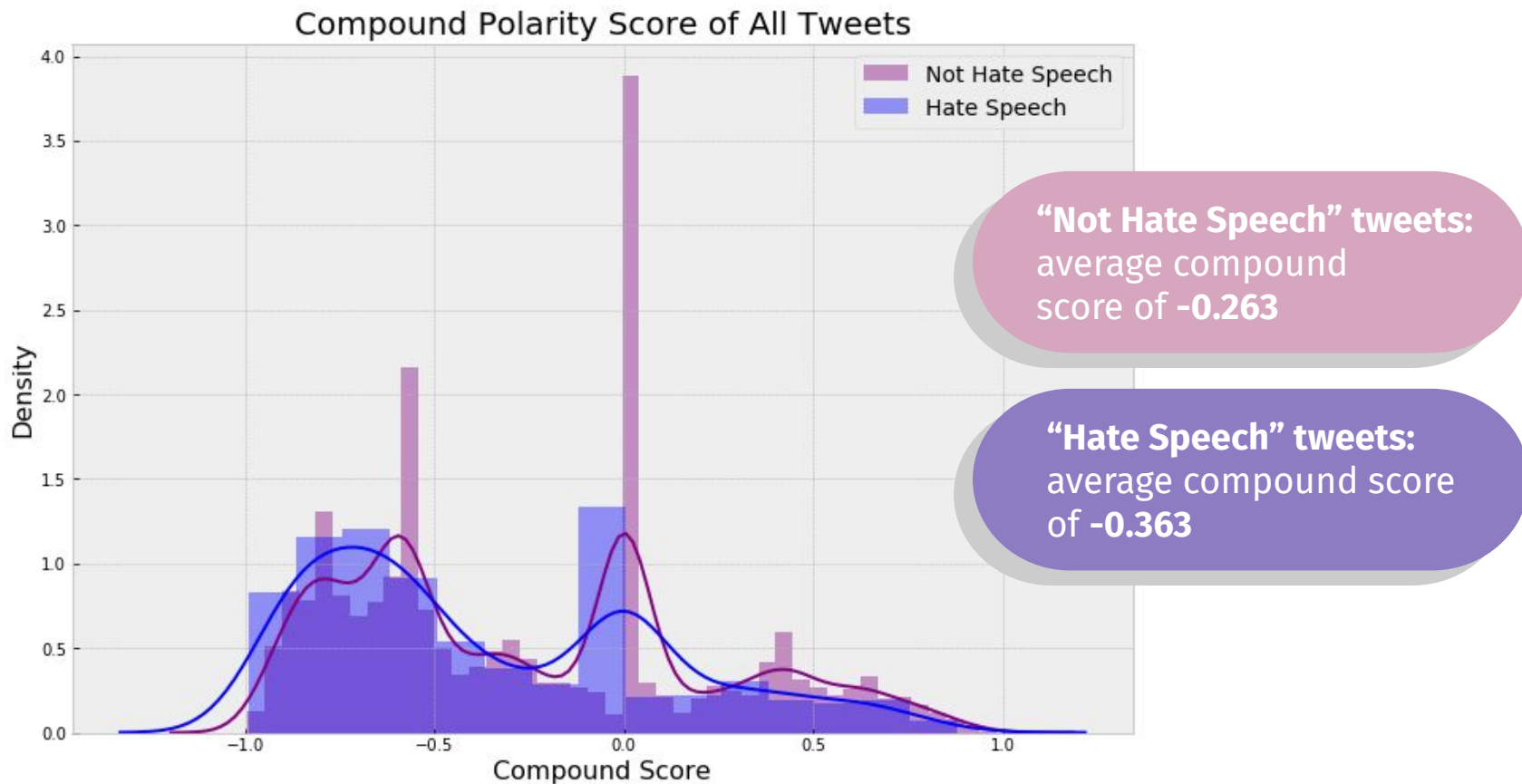
# 1. What are the linguistic differences between hate speech and offensive language?



## 2. What are the most popular hashtags of each tweet type?

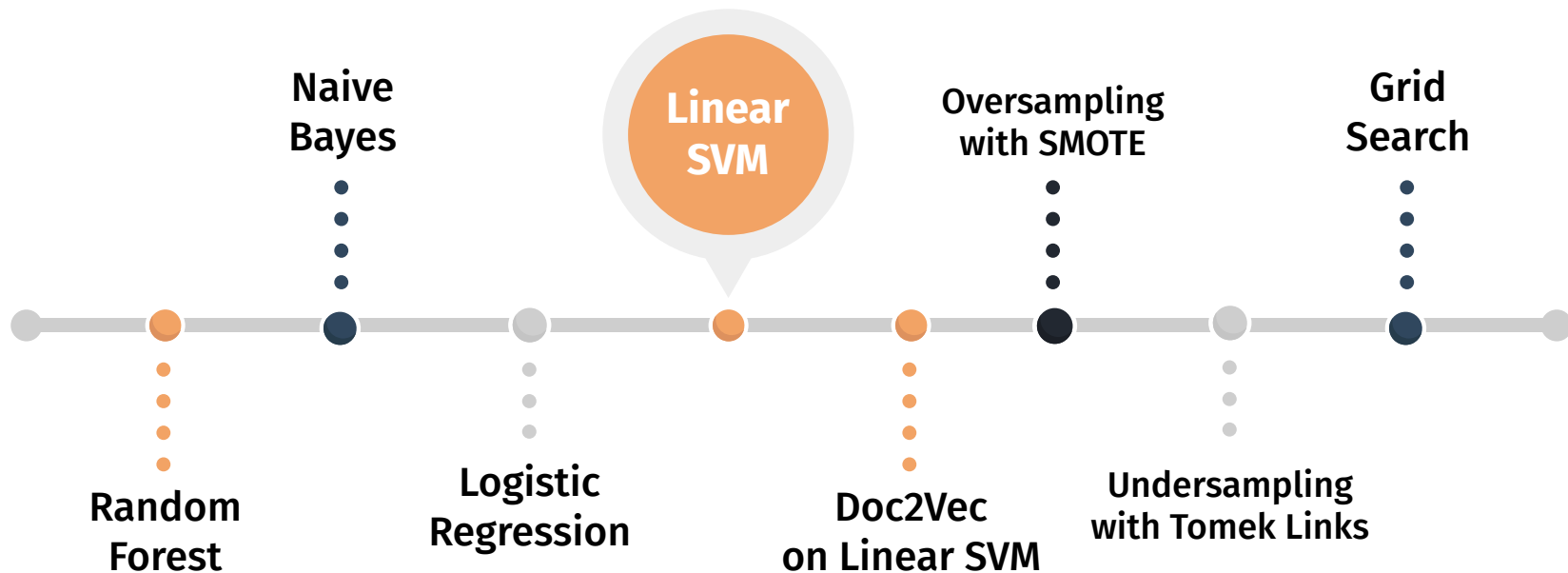


### 3. What is the overall polarity of the tweets?





# Modeling Process



# Final Model Analysis

## Linear SVM Classifier

F1 Score: 0.3955

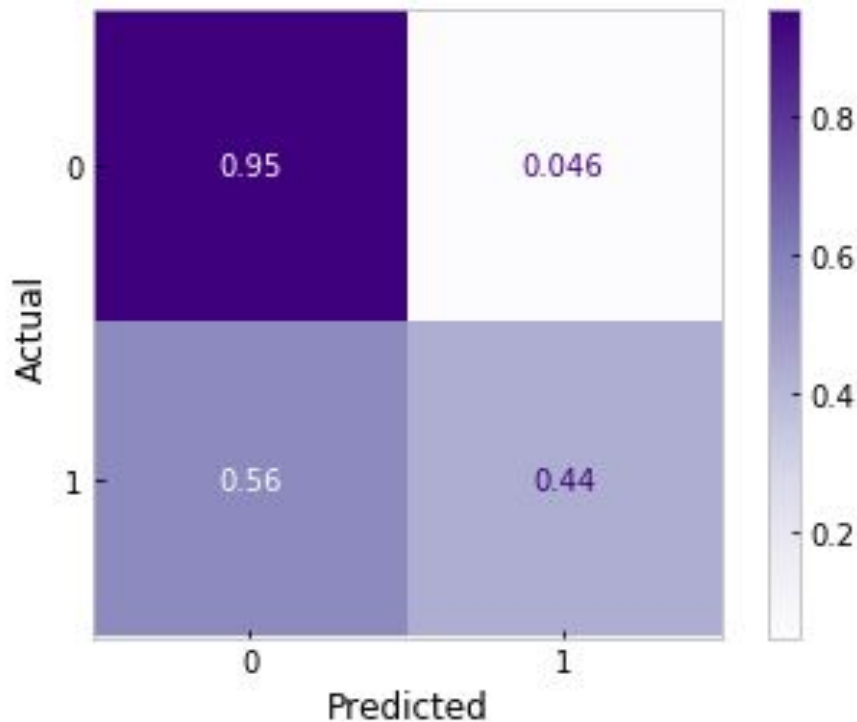
Recall: 0.437

High **True Negative** Rate

Low **False Positive** Rate

Moderate **True Positive** Rate

Normalized Confusion Matrix for Linear SVM

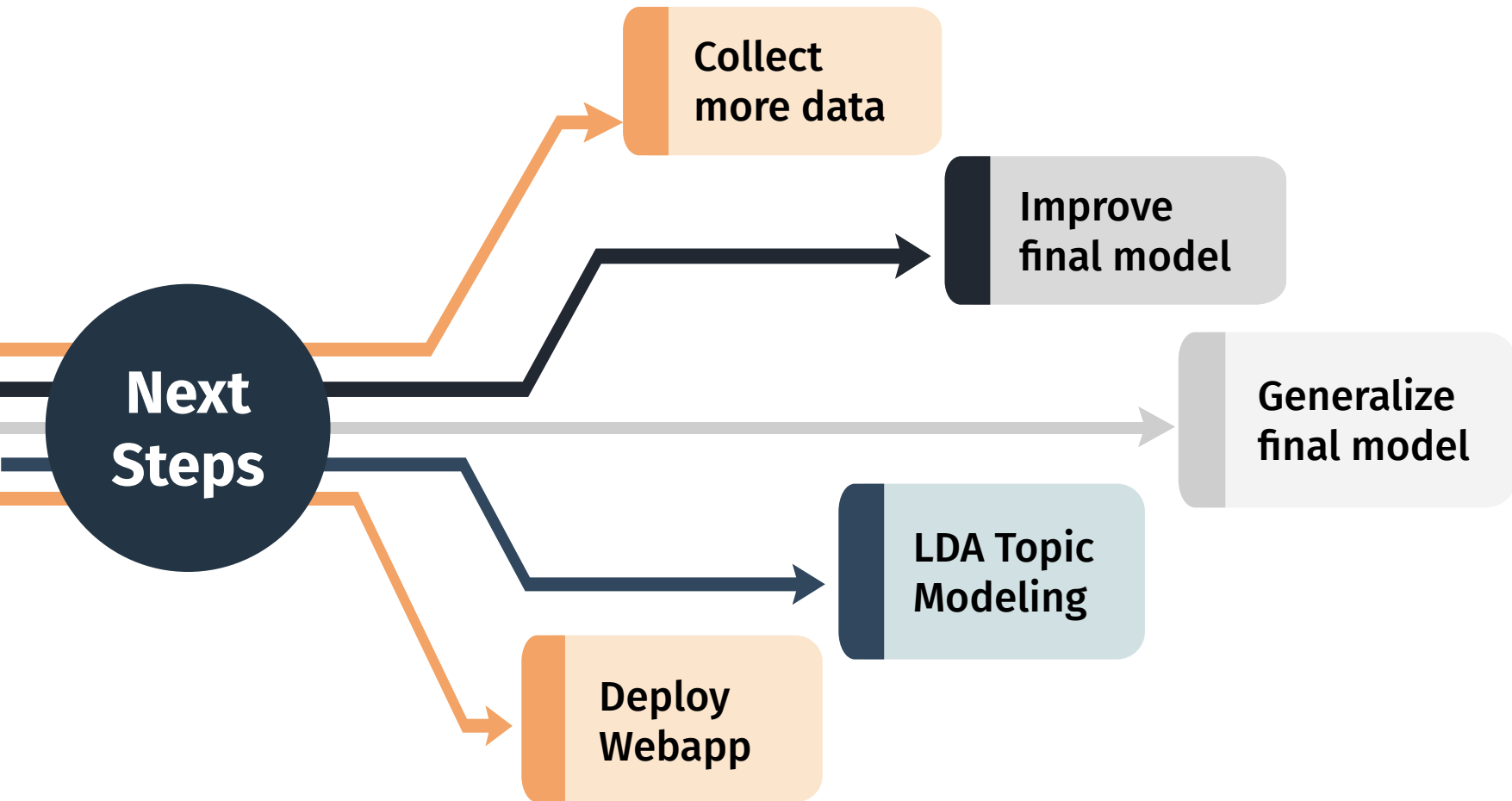


# Conclusion

## Two major roadblocks:

1. Massive class imbalance
2. Model's ability to “understand” hate speech





# Thank You!



**GitHub Repository**

[https://github.com/sidneykung/twitter\\_hate\\_speech\\_detection](https://github.com/sidneykung/twitter_hate_speech_detection)



**SidneyJKung@gmail.com**



<https://www.linkedin.com/in/sidneykung/>



**@Sidney\_K98**

**Presentation Template:  
SlidesGo**