# #AnalyzingHate SpeechOnTwitter

@ChristopherDeLaCruz (not my actual twitter handle)

twitter

# #WhyHateSpeech?

- Over the past few years, multiple organizations (including the FBI) have reported increases in hate speech

- Hate crimes have also increased and generally reflect similar patterns shown in hate speech trends

- In order to do a better job of protecting these classes, we must map a better understanding of current trends and patterns in hate speech

twitter

# #WhatIsHateSpeech?

Is calling someone dumb hate speech?

Depends.

Hate speech is derogatory language aimed at a person based on an identification class (race, gender, class, religion, sexual orientation, etc)

twitter

# #HateSpeechData

Approximately 70,000 labeled tweets (Hate tweets and Not Hate Tweets) were gathered from four different Kaggle datasets and one dataset from Georgia Tech's CLAWS
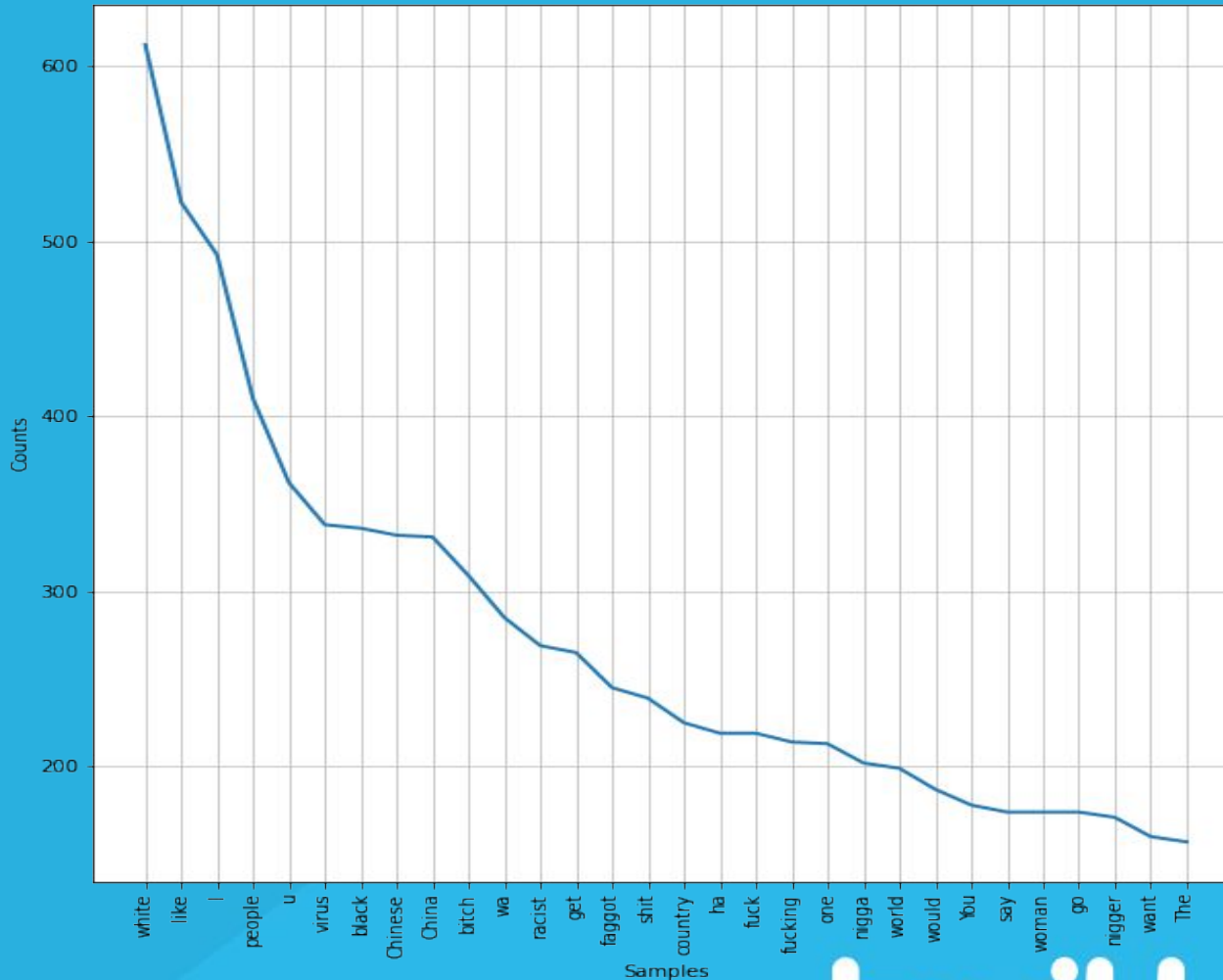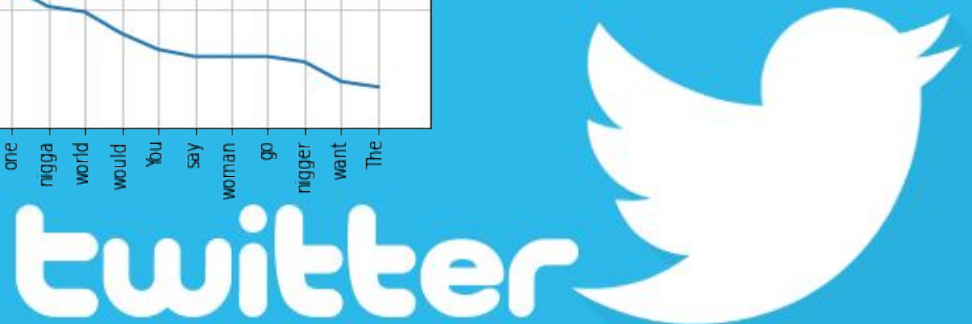
Labels are:
1 - Hate Speech
0 - Not Hate Speech

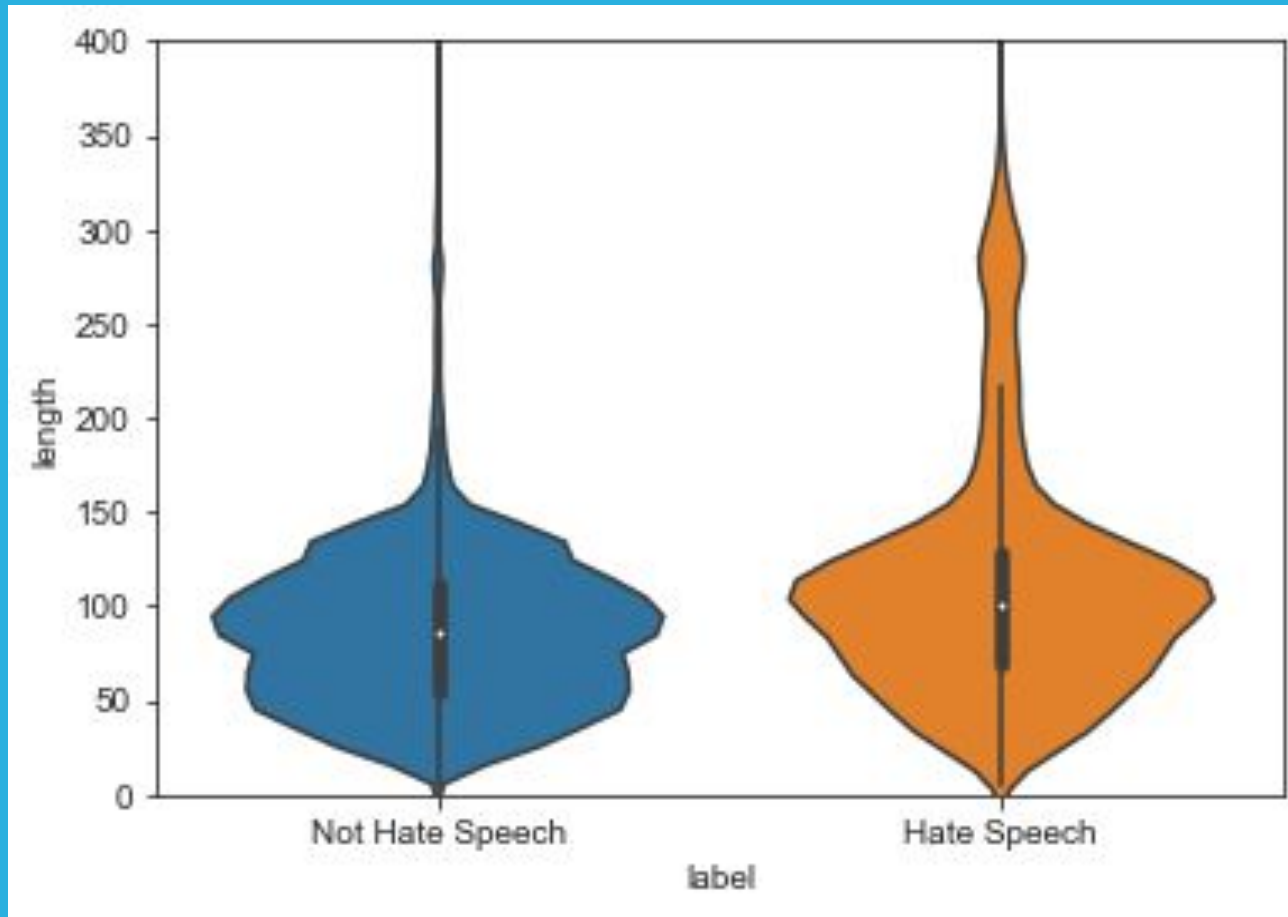twitter

# #HateSpeechResults



## Identity Classes Mentioned Most Often:
- Race: Asian
- Race: Black
- Race: White

Gender: Women
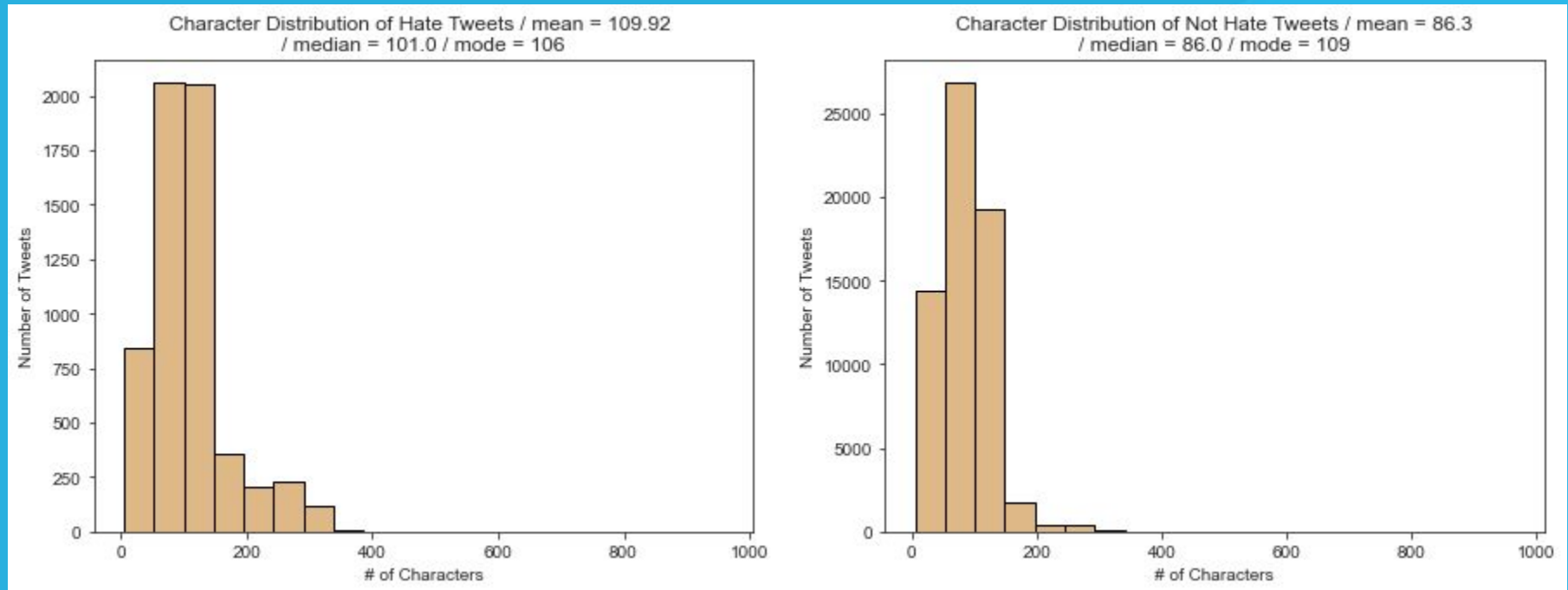- Sexual Orientation: Gay

twitter

# #HateSpeechResults



Character Distribution of Hate Tweets / mean = 109.92 / median = 101.0 / mode = 106

Character Distribution of Not Hate Tweets / mean = 86.3 / median = 86.0 / mode = 109

Hate tweets are about 24 characters longer on average

twitter

# #HateSpeechResults

Over 7,000 models were run but there could only be ONE winner.

First priority - high precision in detecting hate speech

Second priority - high recall in detecting hate speech

Ultimately the best model was...

twitter

# #HateSpeechResults

Support Vector Classifier!

- Able to predict hate speech with 85% precision!

- Recall is 33% and cannot be brought up without sacrificing precision

twitter

# #NextSteps

- Looking into ways to continue improving the model's ability to detect hate speech (more data, different sampling techniques, etc)

- Deploying the model on freshly scraped tweets to gather additional information about hate speech and establishing deeper patterns

twitter