# That's What Who Said?

By Levi Hatch
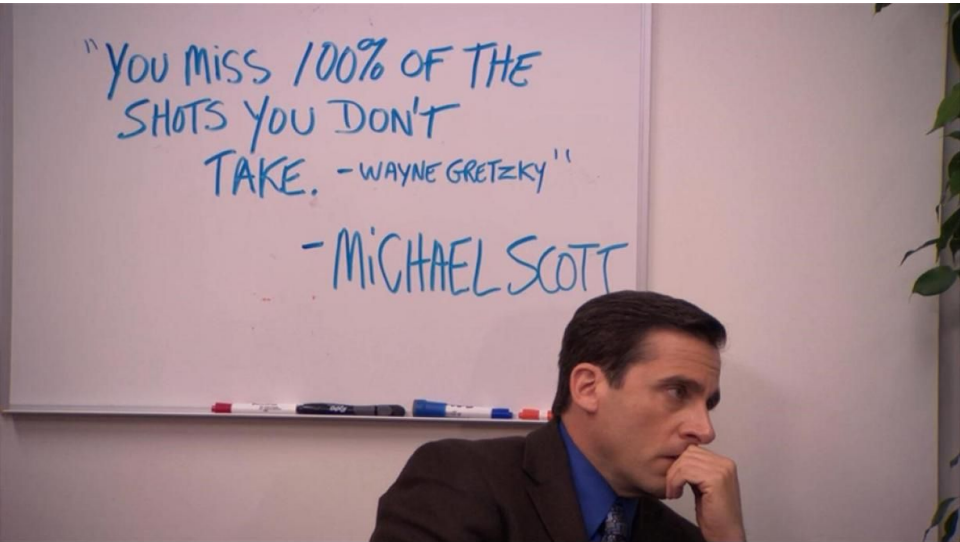
Email: levitravishatch@gmail.com | LinkedIn: /levi-hatch | GitHub: /lthatch

**The inspiration for this project- The Office TV Show**

I am an Office fanatic! I love all the characters and the crazy situations they find themselves in. I would be embarrassed to say how many hours I have watched, but with that time, I feel as if I can quote just about any character from any noteworthy scene.
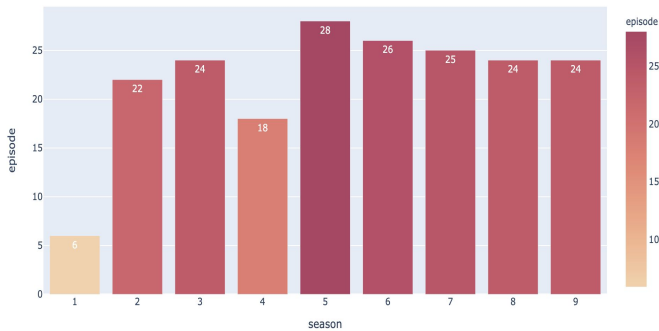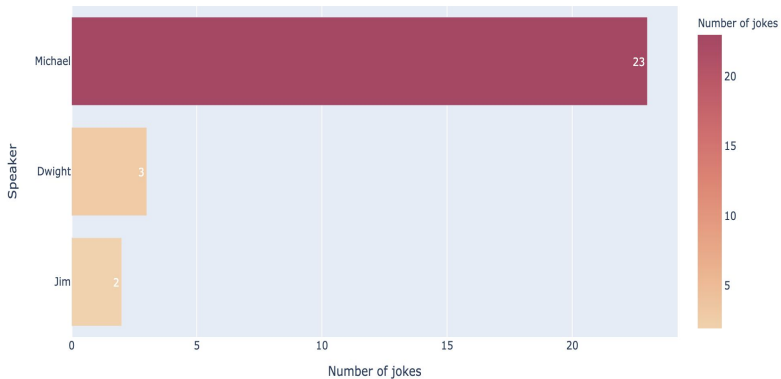
# Project Goal



The goal of this project is to build a model that would take input text from a user and identify a character from The Office who would most likely say those words.
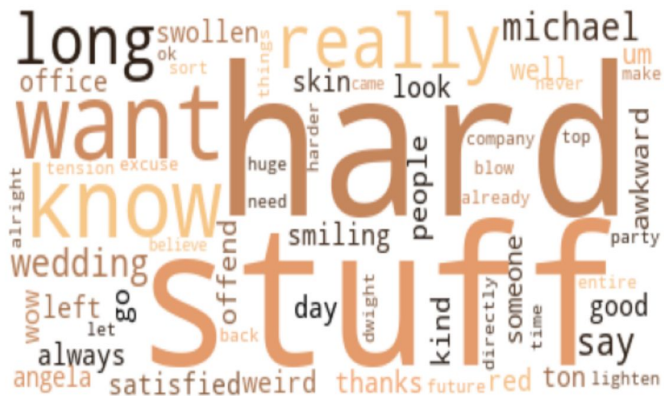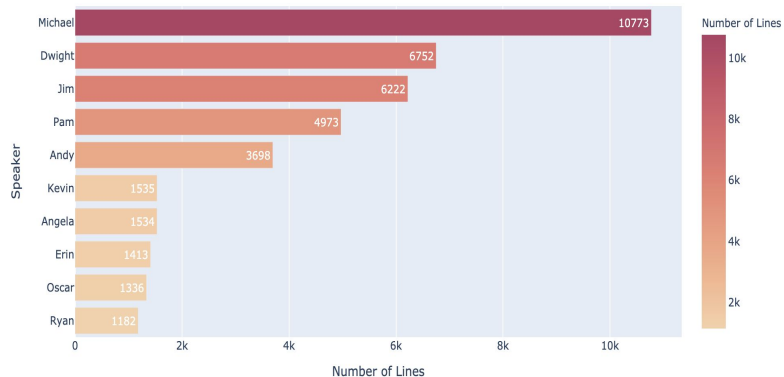
# Exploratory Data Analysis

# Preparing the Data



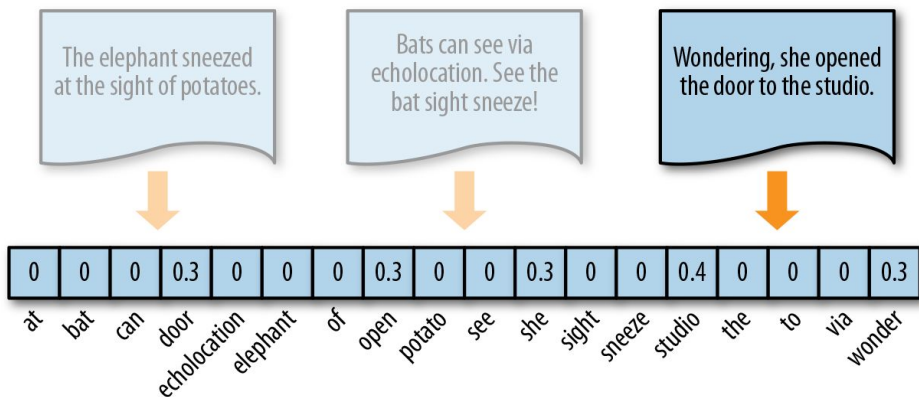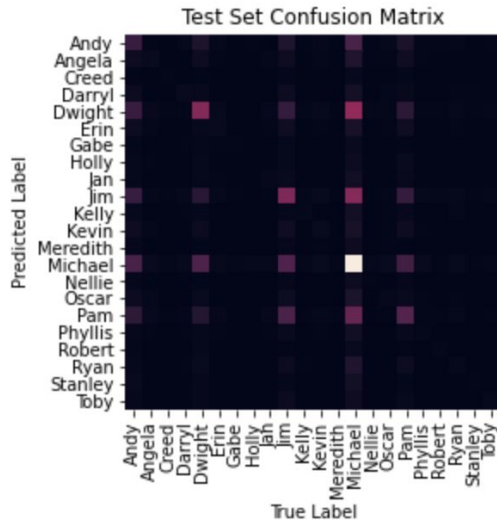- I limited the amount of classes/characters down to the top 20 speakers
- Using a TfidfVectorizor, I removed stop words and transformed the text into a weighted matrix of numbers
- Once in numeric form, I performed a train, test split and my data was ready to be fed into a model
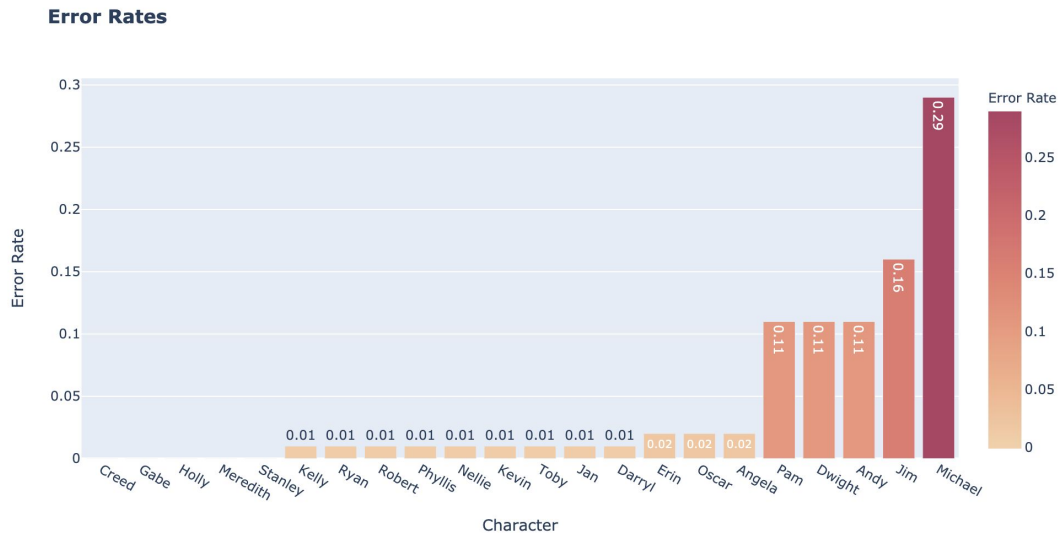
# Selecting and Training a Model



Test Set Confusion Matrix

Hold out data accuracy: **.24**

Random guessing expected accuracy: **.05**

I decided to use scikit learn's ComplementNB model because of the imbalances in my data.



**Error Rates**

# Limiting the Number of Classes


Test Set Confusion Matrix

Hold out data accuracy: **.41**

Random guessing expected accuracy: **.25**

In an effort to improve my accuracy I decided to limit the number of classes the model needed to predict. I chose the top 4 speakers


Error Rates

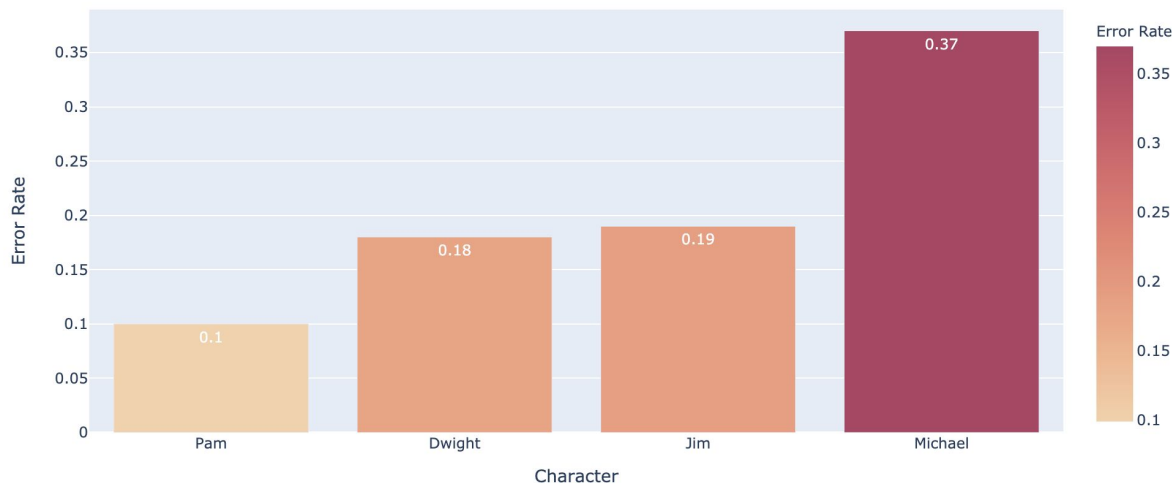# Fine-tuning the model

After seeing my accuracy increase I decided to add in 3 more characters. In addition to increasing the number of classifications, I tuned the model using sklearn's Randomized SearchCV
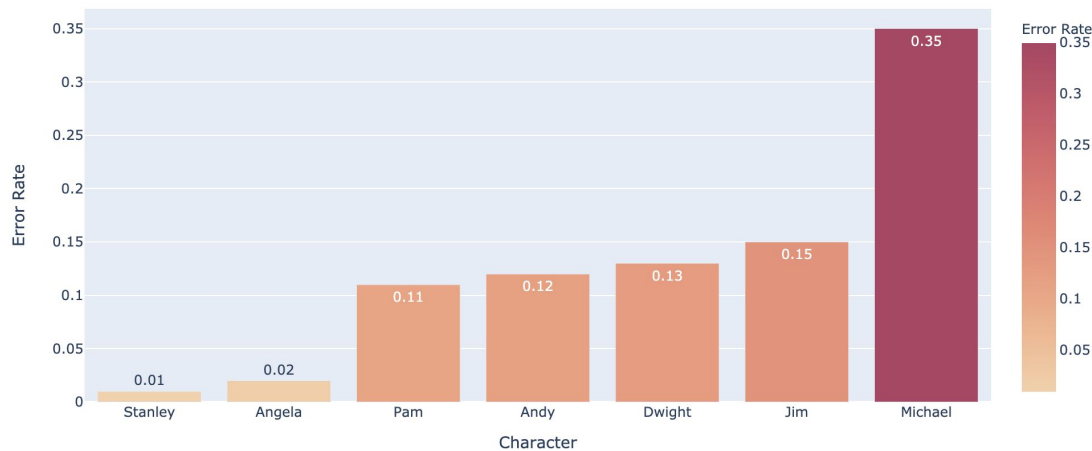


Test Set Confusion Matrix

Hold out data accuracy: **.33**

Random guessing expected accuracy: **.14**



Error Rates

# In Conclusion of this Powerpoint

As you can see, the results from the flask app were mostly as expected.

**Next steps:** Expanding this project into creating a chat bot in order to chat with your favorite characters!



Email: levitravishatch@gmail.com | LinkedIn: /levi-hatch | GitHub: /lthatch

[Flask App](Flask App)