

THE ONE ABOUT F•R•I•E•N•D•S

Character Prediction
based on Dialogue





1. INTRODUCTION
2. DATA ACQUISITION



3. DATA CLEANING



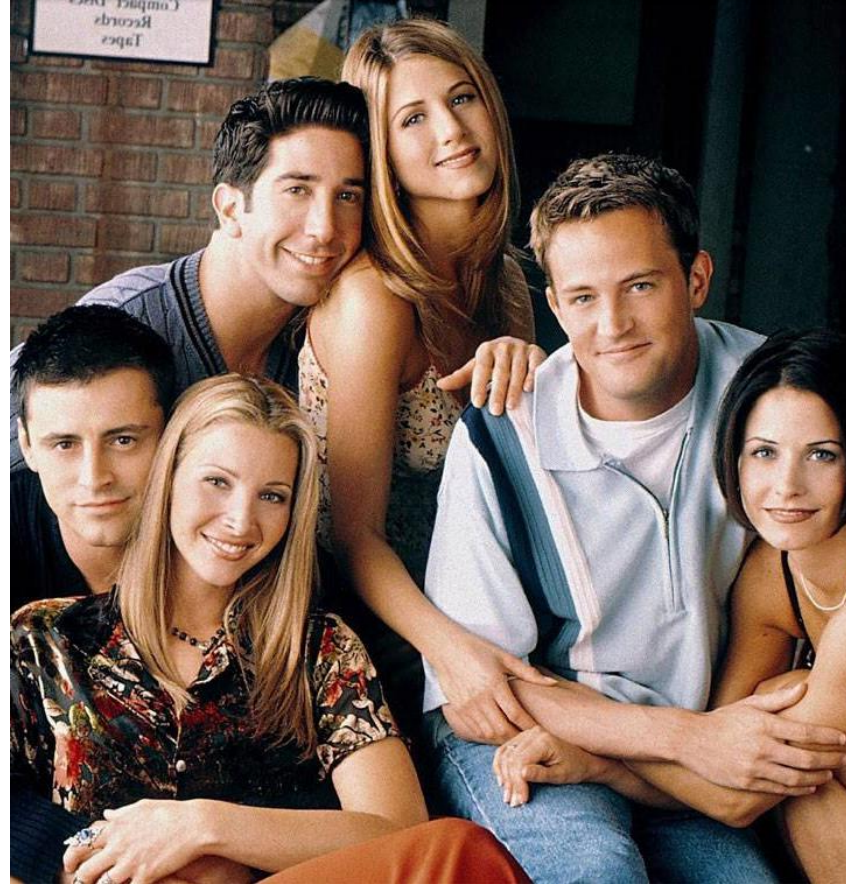
4. MACHINE LEARNING

*Oh My God!
Seriously?!*



5. RESULTS
6. WHAT'S NEXT?

THE ● ONE ● WITH
THE ● INTRODUCTION



ARE YOU A FRIENDS FANATIC?

THE ● ONE ● WITH ● ALL ● THE ● LINES

- 236 Episodes
- Combination of BeautifulSoup and Regex
- Several episodes were “irregular”
- Used aggregation functions to determine if I needed to re-scrape any data. Total episodes not included = 6, which was only 3% of the total

```
Characters = []
Lines = []
Title = []
Season = []
Episode = []
for url in tqdm(all_episodes_urls):
    r = requests.get(url)
    soup = BeautifulSoup(r.text, 'html.parser')
    ps = soup.find_all('p')
    for p in ps:
        char = regex.findall(r"[A-Z][a-zA-Z. ]+:", p.text)
        if char != []:
            if char[0] != "Scene:":
                Characters.append(char[0])
                index = regex.search(char[0], p.text).start() + len(char[0])
                line = p.text[index:]
                Lines.append(line.replace("\n", " "))
                Title.append(soup.title.string)
                season = regex.findall('friends/(\w+\d+)', url)
                Season.append(season)
                ep = regex.findall('friends/\w+\d+/(\d+)', url)
                Episode.append(ep)
```

THE ● ONE ● WITH ● ALL ● THE ● CLEANING

- Reviewed the data using methods noted below
- Created a Name function to correct misspellings
- Used a series of `regex.sub(pattern, "", line)` to replace non-dialogue

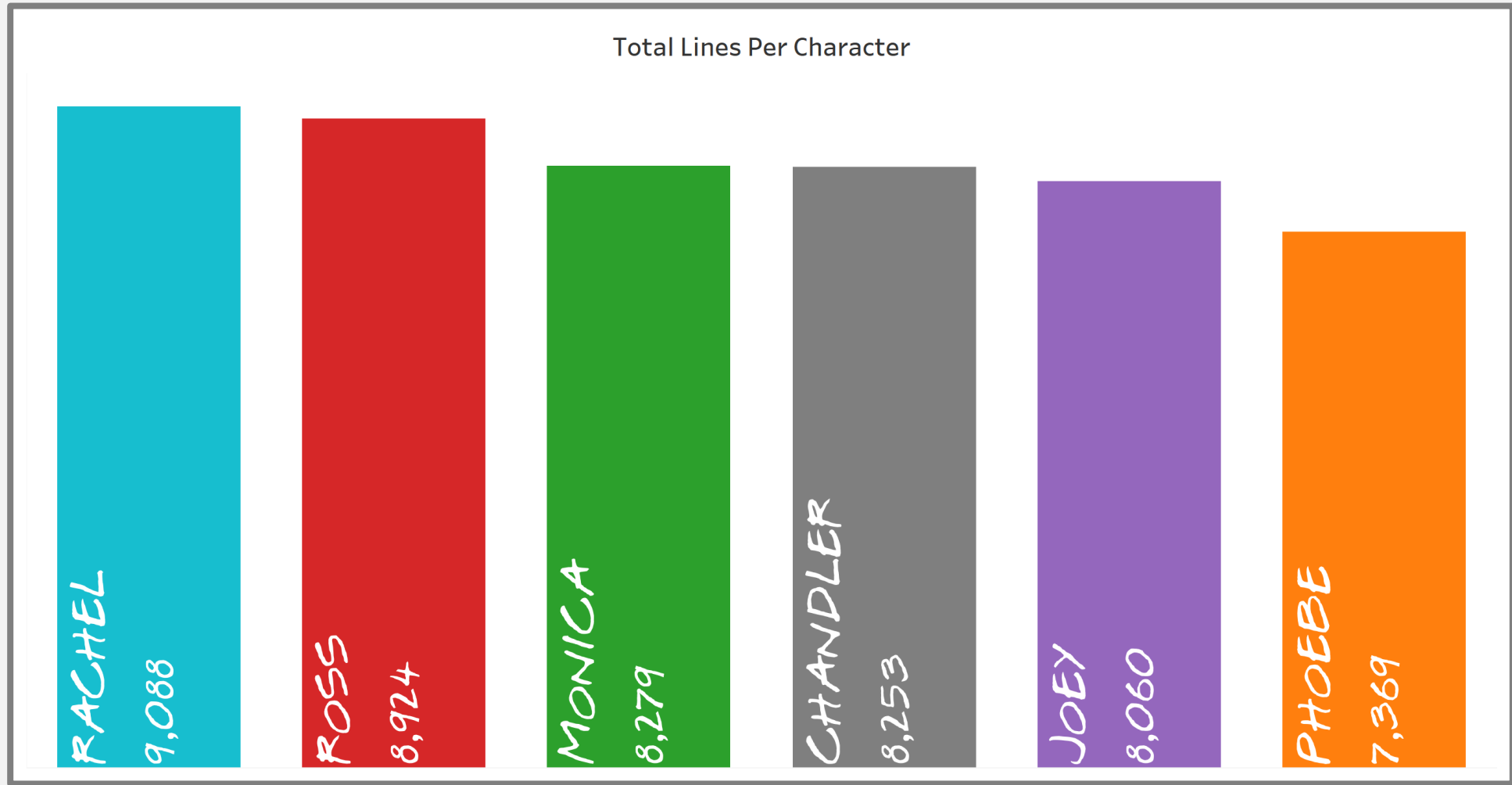
- Dropped all null lines and lines from "Other"

cast members

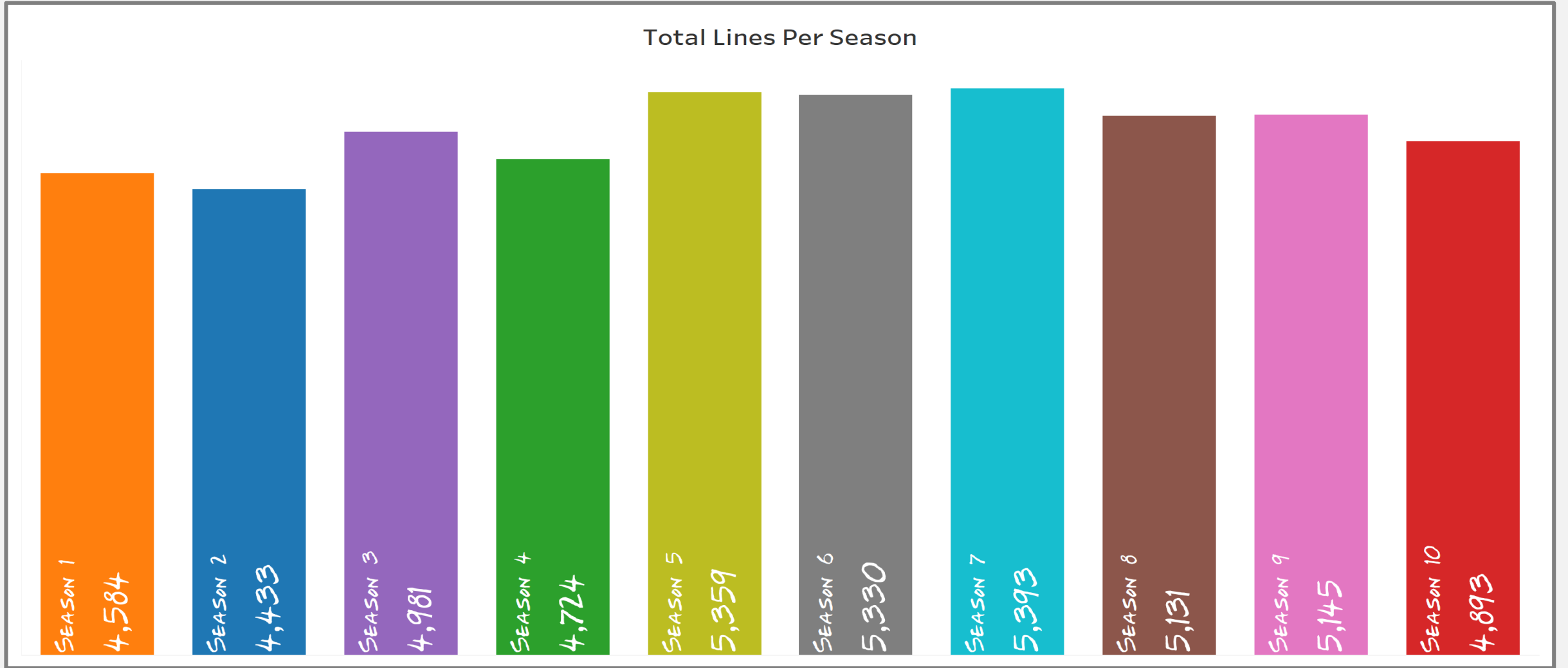
- `df.info()`
- `df.isnull().sum()`
- `df.dropna()`
- `df.Character.isin(main_characters).sum()`
- `set(df.Character)`
- `def name_func(char)`
- `df.groupby('Character').count().loc[main_characters]`

```
Ok. (She sits on the bed and Ross sits near her)<+i>
Thank you for coming with me today.
(to Ross) Let me get you some coffee.
About 20 minutes. <font size="4"> <b> CLOSING CREDI
TS </b> </font>
Well, I tracked down Marcel and get this, he's healt
hy, he's happy, and he's right there in New York fil
ming <i> Outbreak II - The Virus Takes Manhattan </i>
> .
```

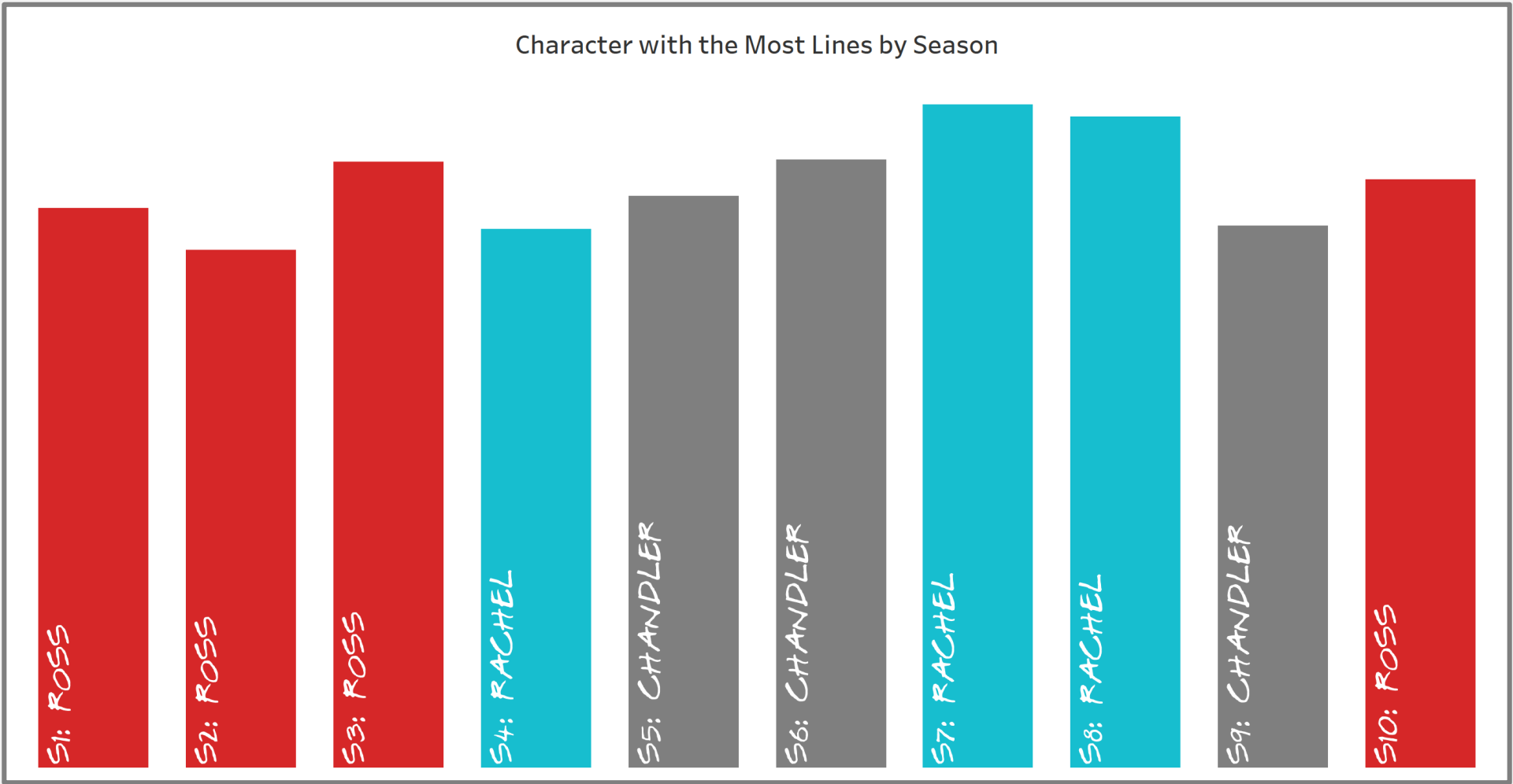
THE • ONE • WITH • THE • EXPLORATION • MISSION



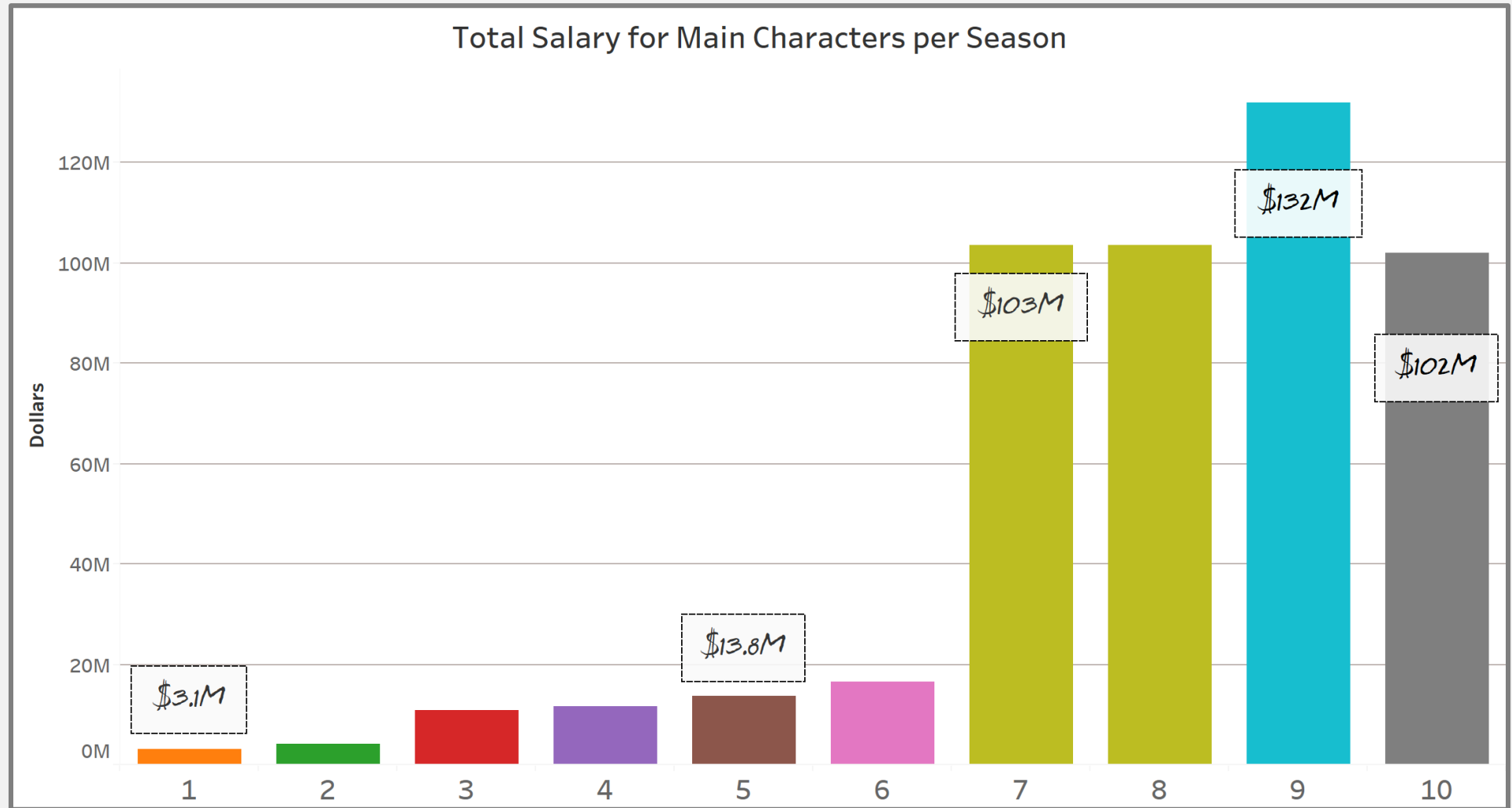
THE ● ONE ● WITH ● THE ● EXPLORATION ● MISSION



THE ● ONE ● WITH ● THE ● EXPLORATION ● MISSION



THE ● ONE ● WITH ● THE ● EXPLORATION ● MISSION



THE ● ONE ● WITH ● THE ● MODELS

	Character	Baseline Percent
0	Rachel	18.2%
1	Ross	17.9%
2	Monica	16.6%
3	Chandler	16.5%
4	Joey	16.1%
5	Phoebe	14.7%
6	Average	16.7%

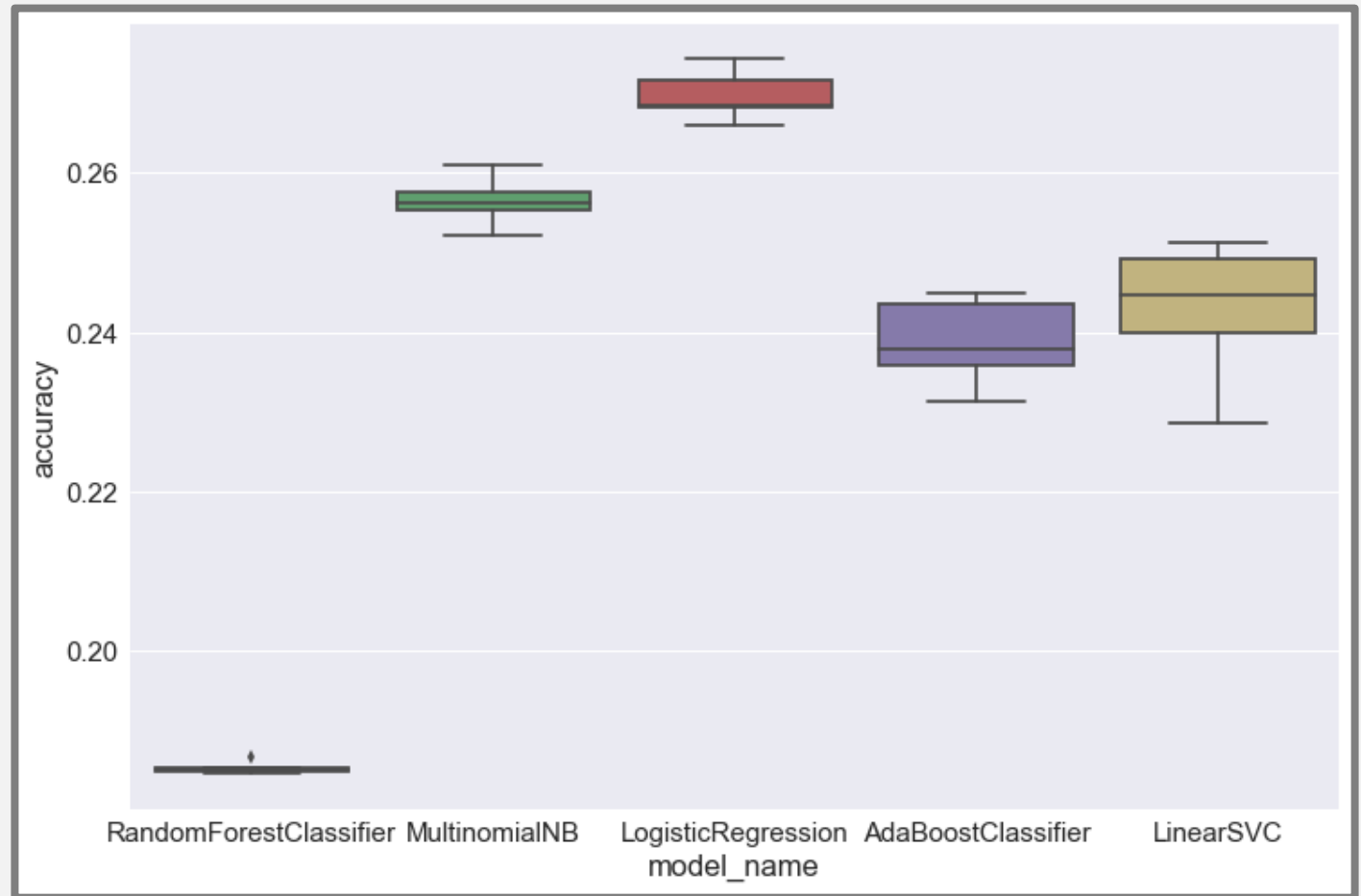
Models Used

- Logistic Regression
- Decision Tree Classifier
- Linear SVC
- Random Forest Classifier
- Extra Trees Classifier
- AdaBoost
- Multinomial Naive Bayes
- SGD Classifier

THE ONE WITH THE MODELS

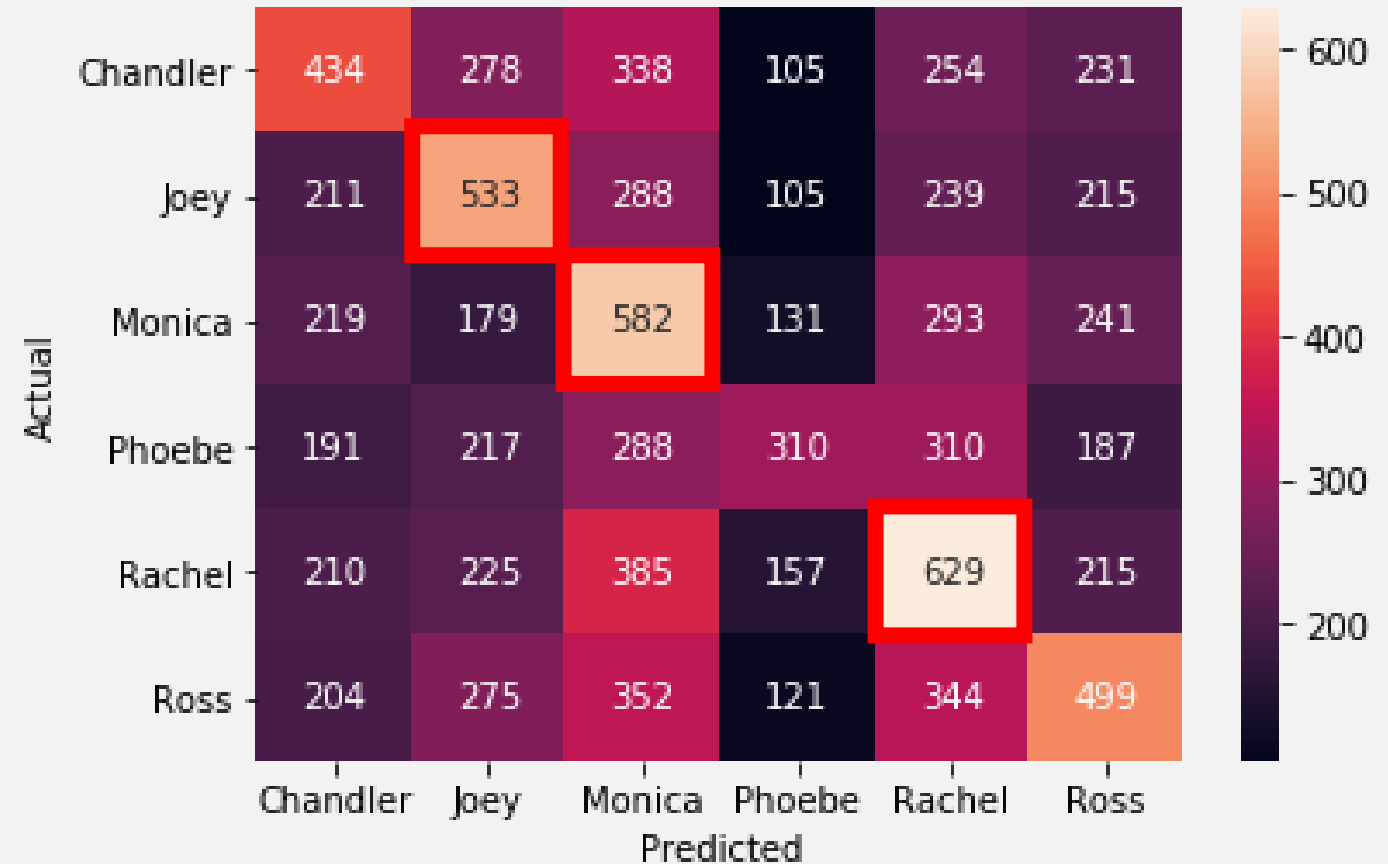
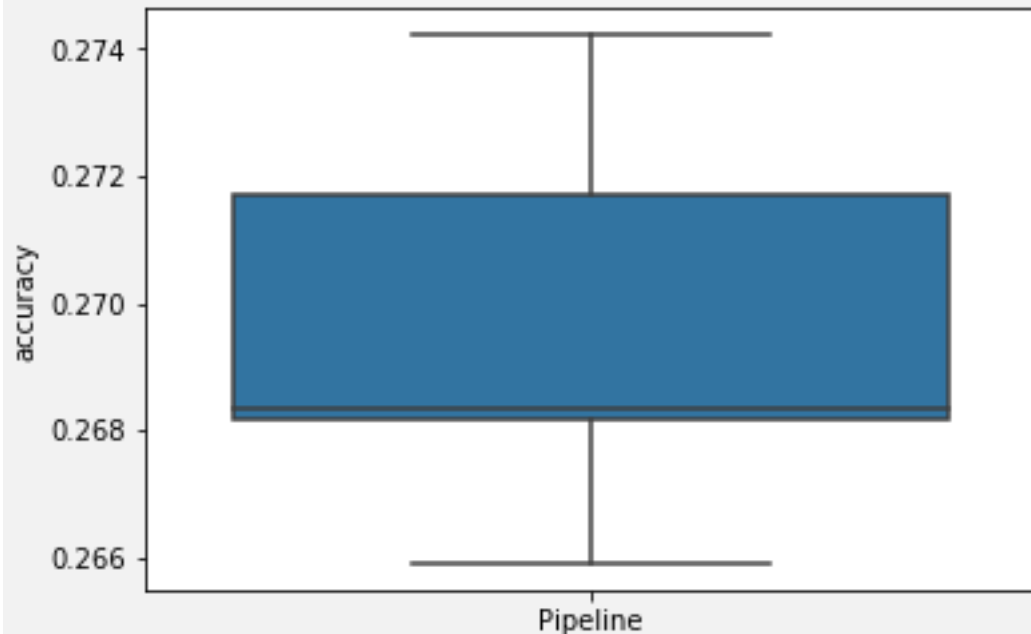
What I learned:

1. Gridsearch is **slow**
2. Pipelines are more efficient
3. PCA – lose information, didn't use to test predictions
4. Know when to **STOP**



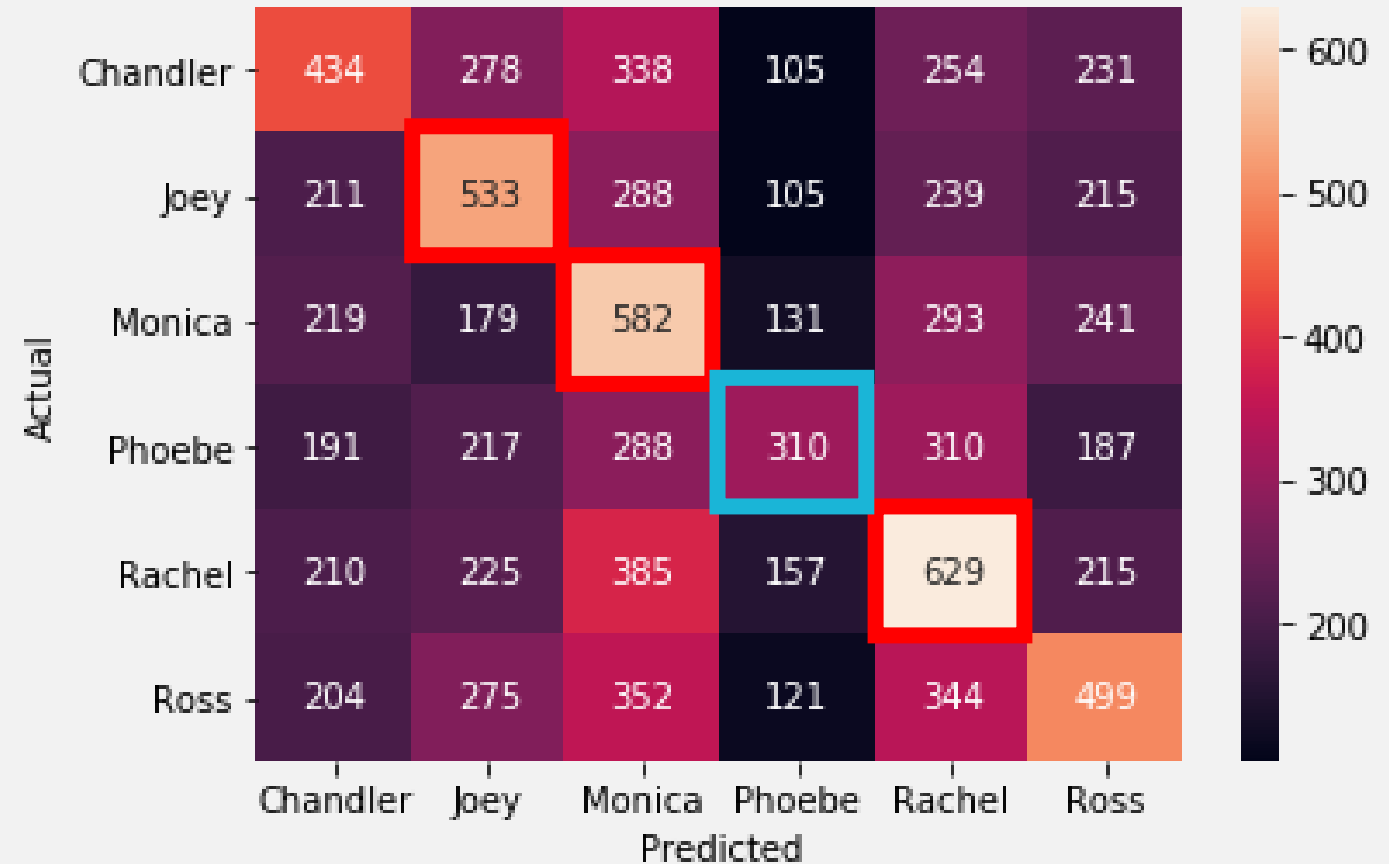
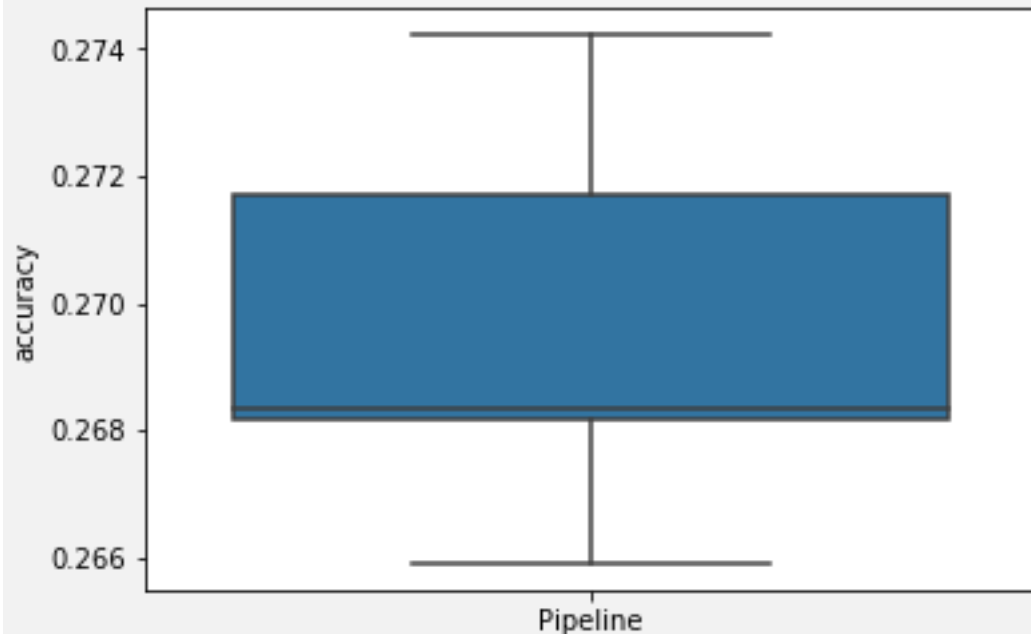
THE ONE WITH THE RESULTS

- The Winner is.... Logistic Regression!
 - Preprocessor: Count Vectorizer
 - Ngram Range: 1-3
 - Stop Words = English
 - Accuracy Score: 0.2988
 - Number of features: 219,010



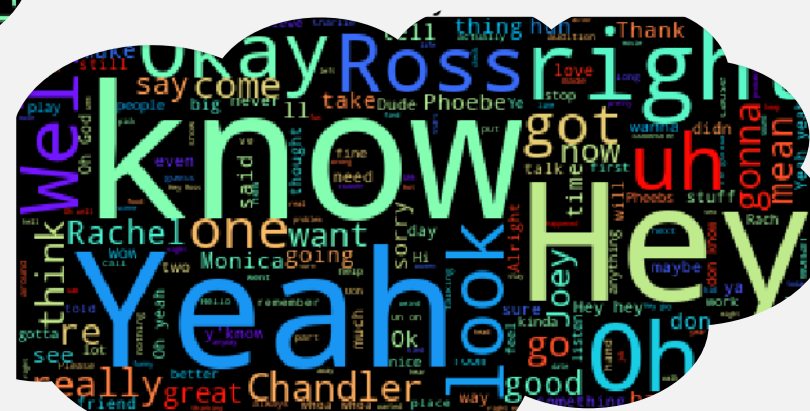
THE ONE WITH THE RESULTS

- The Winner is.... Logistic Regression!
 - Preprocessor: Count Vectorizer
 - Ngram Range: 1-3
 - Stop Words = English
 - Accuracy Score: 0.2988
 - Number of features: 219,010



WORD CLOUDS

Can you tell who is who?



Phoebe



Chandler



Rachel



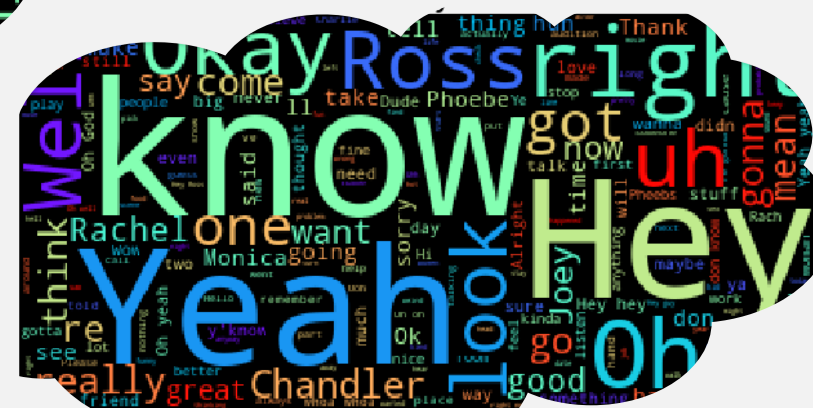
Monica



Ross



Joey



FLASK•ON•FLEEK

Friends Me!

Enter Your Phrase Below

Insert Phrase

Friend Me!

Enter Phrase

INPUT PAGE

Friends Me!

Character: Chandler

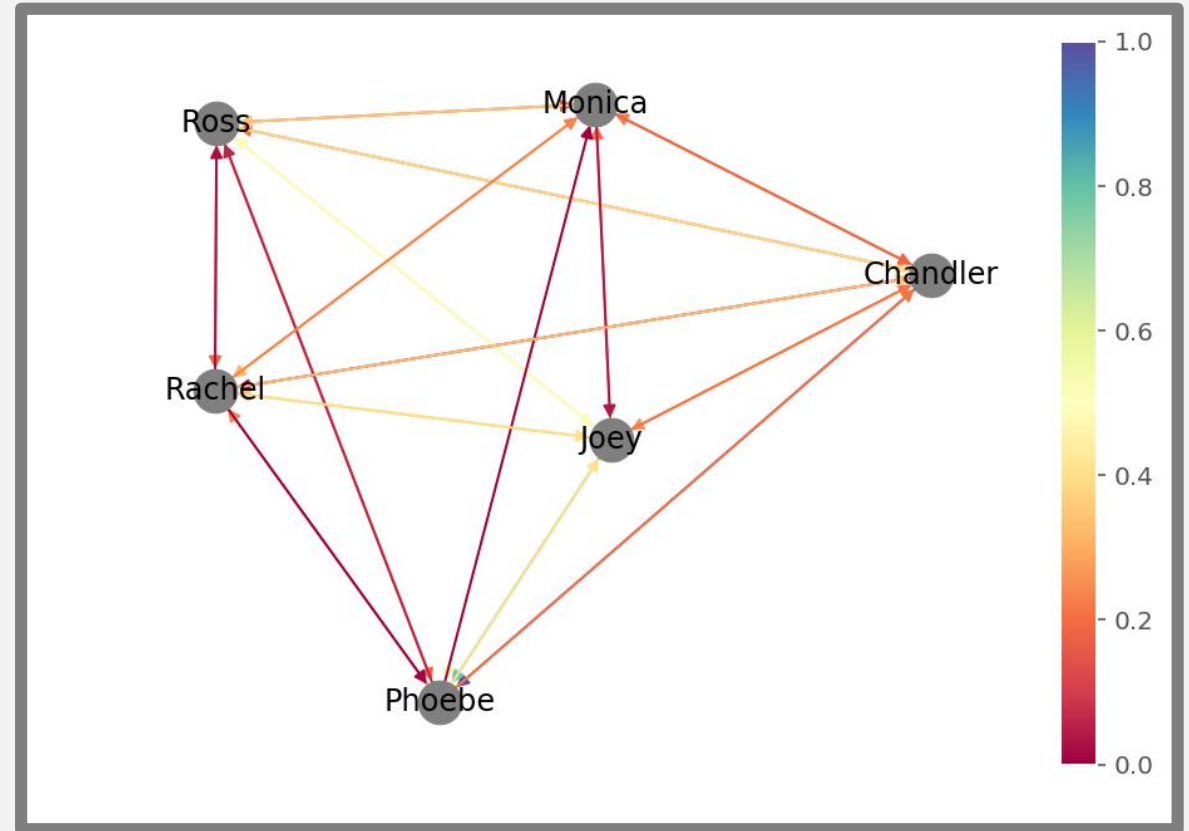


Get a Prediction!

PREDICTION PAGE

THE • ONE • WITH • THE • FUTURE... STEPS

- What else can we do?
 - Download someone's twitter feed to predict which Friends character they are most like
 - Determine the centrality of the characters



THANKS!

