

# Using NLP to predict Yelp Stars Given

FlatIron School- Capstone Project



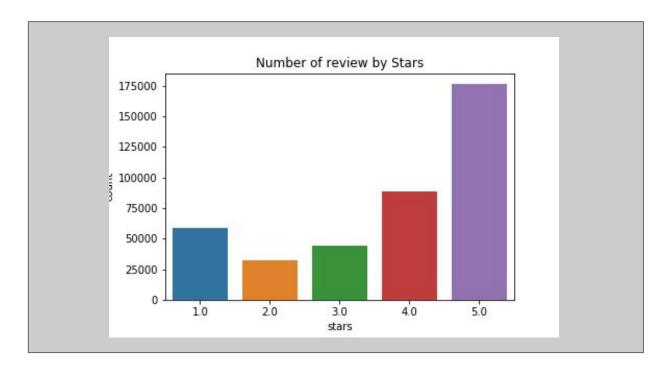
### Purpose and Process

Can we predict number of stars given through Natural Language Processing?



- What is the breakdown of reviews given by customers?
- Is there a difference in language used based on the number of stars given?
- Are there any important word associations and how do we capture those?
- Are there any themes that emerge from the data and do those themes vary by number of stars?

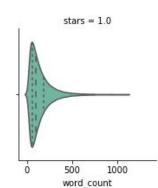


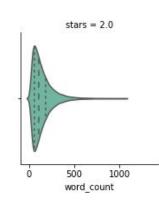


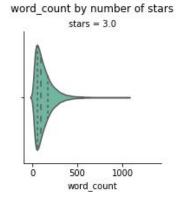


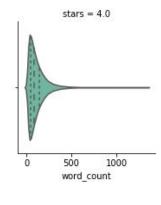
- What is the breakdown of reviews given by customers?
- Is there a difference in language used based on the number of stars given?
- Are there any important word associations and how do we capture those?
- Are there any themes that emerge from the data and do those themes vary by number of stars?

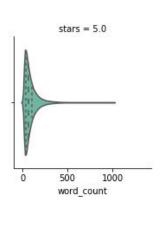


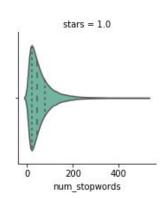


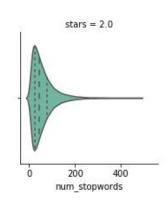


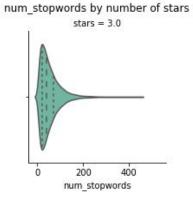


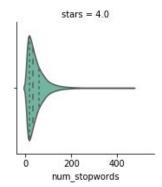


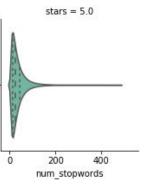




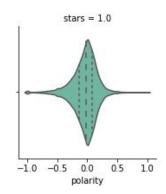


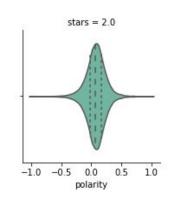


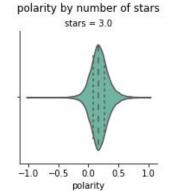


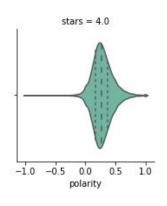


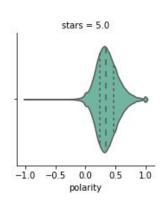


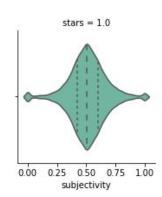


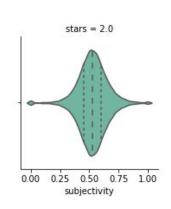


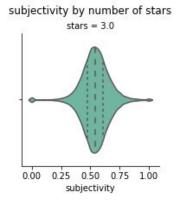


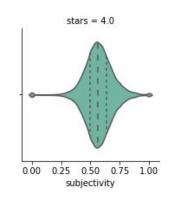


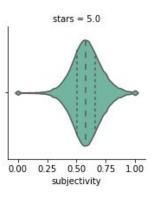














- What is the breakdown of reviews given by customers?
- Is there a difference in language used based on the number of stars given?
- Are there any important word associations and how do we capture those?
- Are there any themes that emerge from the data and do those themes vary by number of stars?



#### Most Similar Words to sandwhich ignoring None

-----

panini: 0.8566

club sandwich: 0.8416 turkey sandwich: 0.8385

tuna salad: 0.8136 cobb salad: 0.8114

roast beef sandwich: 0.8080

blt: 0.8065

pastrami sandwich: 0.8028

turkey club: 0.7944

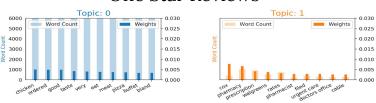
hoagie: 0.7884





- What is the breakdown of reviews given by customers?
- Is there a difference in language used based on the number of stars given?
- Are there any important word associations and how do we capture those?
- Are there any themes that emerge from the data and do those themes vary by number of stars?

#### One Star Reviews



0.025

0.020

0.015

0.010

0.005

0.030

0.025

0.020

0.015

0.010

0.005

0.030

0.025

0.020

0.015

0.010

0.005

0.030

0.025

0.020

0.015

0.010

0.005

Weights

Weights

Topic: 2

Topic: 4

Topic: 6

Topic: 8

do told your their her alled said

Weights

Word Count

Word Count

Word Count

Word Count

5000

4000

3000

1000

6000

5000

4000

3000

2000

1000

6000

5000

4000

3000

2000

1000

6000

5000

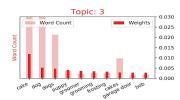
4000

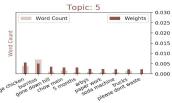
3000

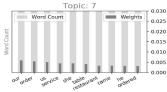
2000

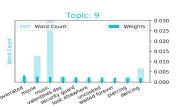
1000

€ 2000



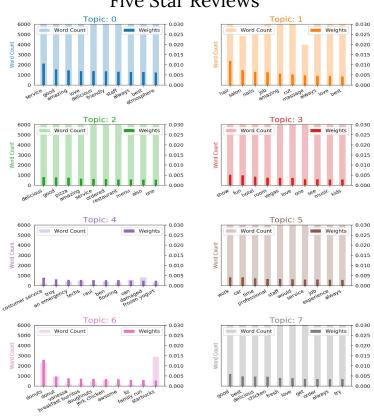


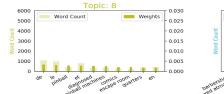


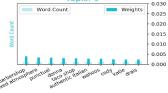




#### Five Star Reviews







Topic: 9

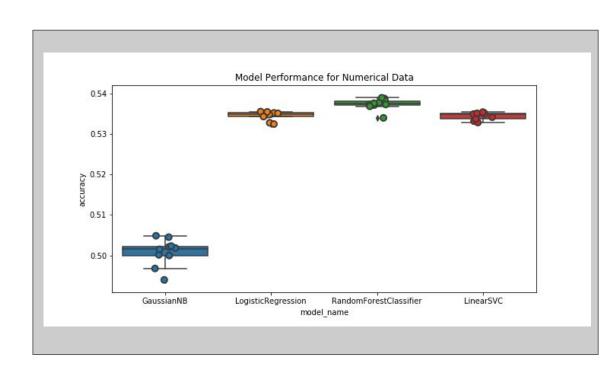


Naive Bayes Model: 50.08% (+/- 0.3167%)

Logistical Regression Model: 53.46% (+/- 0.1023%)

Random Forest Model: 53.74% (+/- 0.1313%)

Linear SVC Model: 53.45% (+/- 0.08733%)



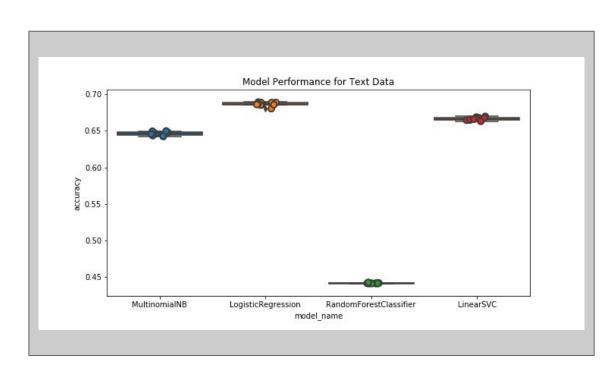


Naive Bayes Model: 64.62% (+/- 0.2514%)

Logistical Regression Model: 68.65% (+/- 0.2543%)

Random Forest Model: 44.19% (+/- 0.04548%)

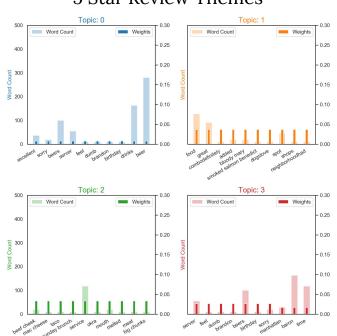
Linear SVC Model: 66.64% (+/- 0.1989%)





## Testing with New Data

#### 5 Star Review Themes







- 1. Work to improve our cleaner function to improve our data.
- 2. Test an unsupervised model that combines both numerical and text data



#### **Any questions?**