



Yelping Around

Group 2

Ishan Shah, Svetak Sundhar, Katherine Wang,
Jessica Ho, Alex Kan



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

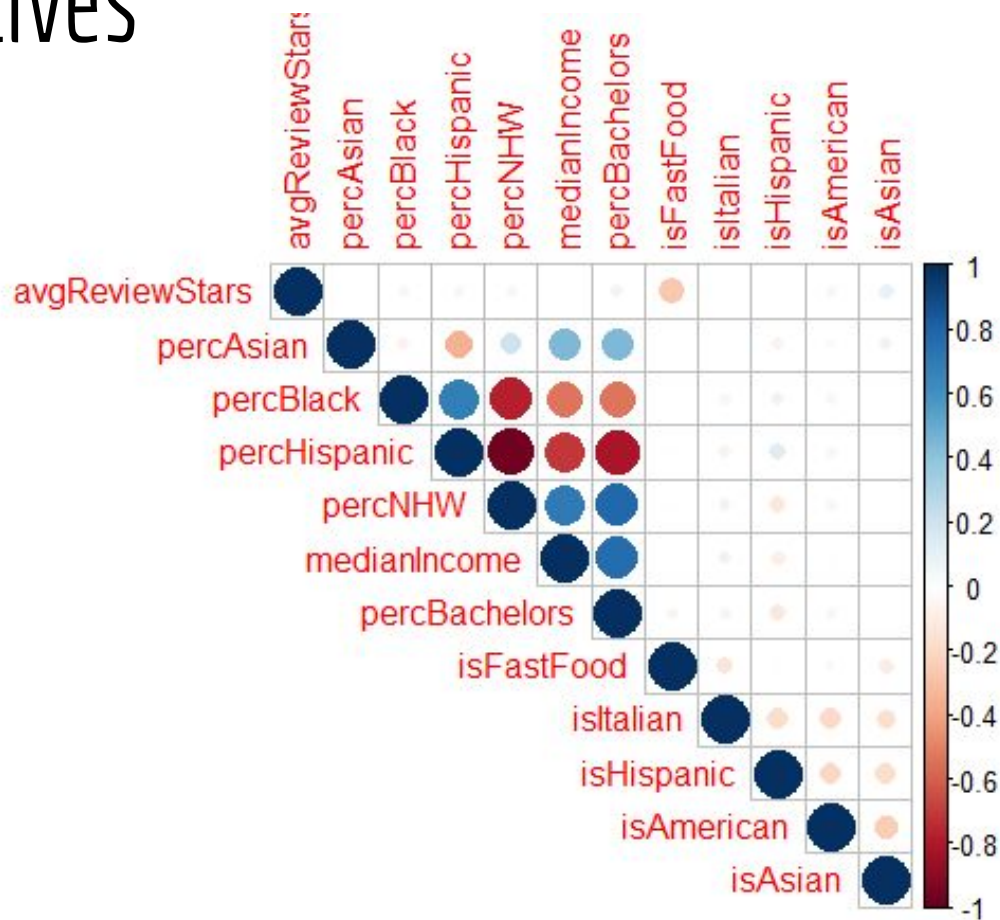
Yelp?

- Tons of places to build a restaurant – why not use Yelp to help?
- What attributes would you need to be successful in a specific geographic location?
 - What demographics should you target?
 - How would you expect to be rated by customers?
 - How could you predict your number and quality of review??

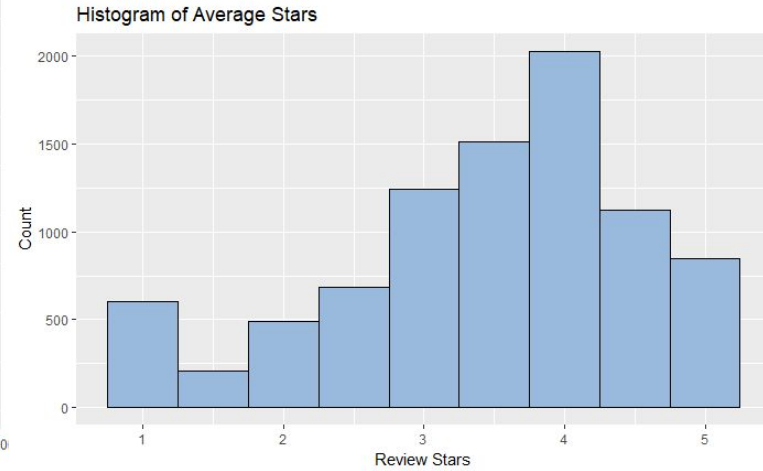
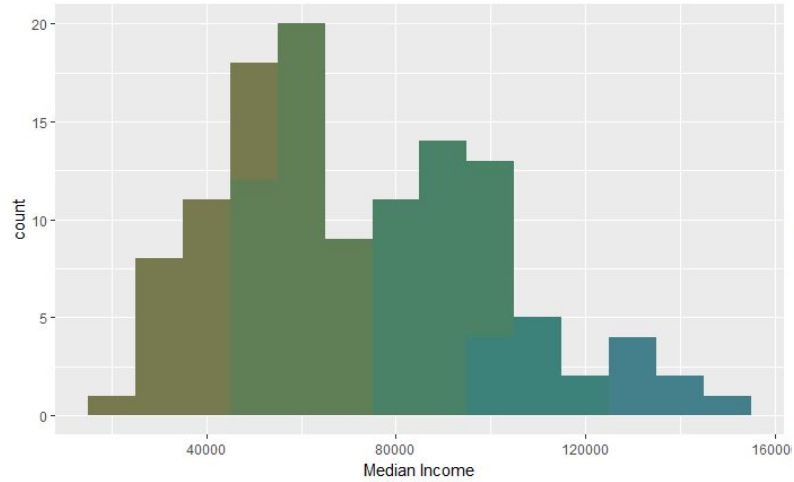
The Data

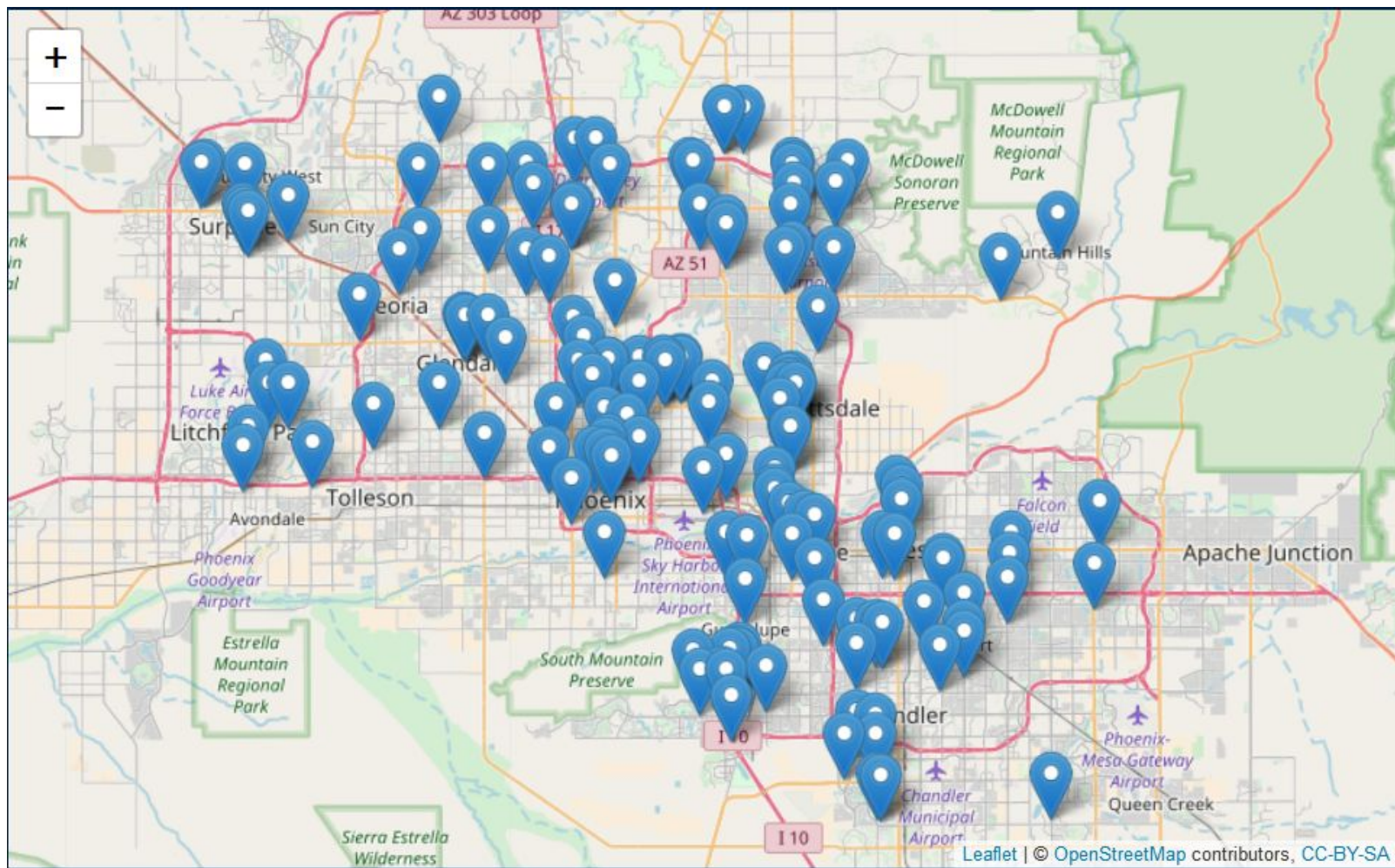
- Yelp Challenge Dataset
 - Tables for all businesses, users, and reviews
- Census Demographic Data
 - Racial distribution, median income, age distribution, and education level by ZIP Code
- Pre-Processing
 - Filtered out restaurants within a 30 mile radius of Phoenix downtown
 - Aggregated by business ID
 - Focused on average star review, average Yelp age of reviewers, characteristics of the reviews, and dummy variables for types of food served
 - Combined both datatables – 8,700 rows, 49 variables

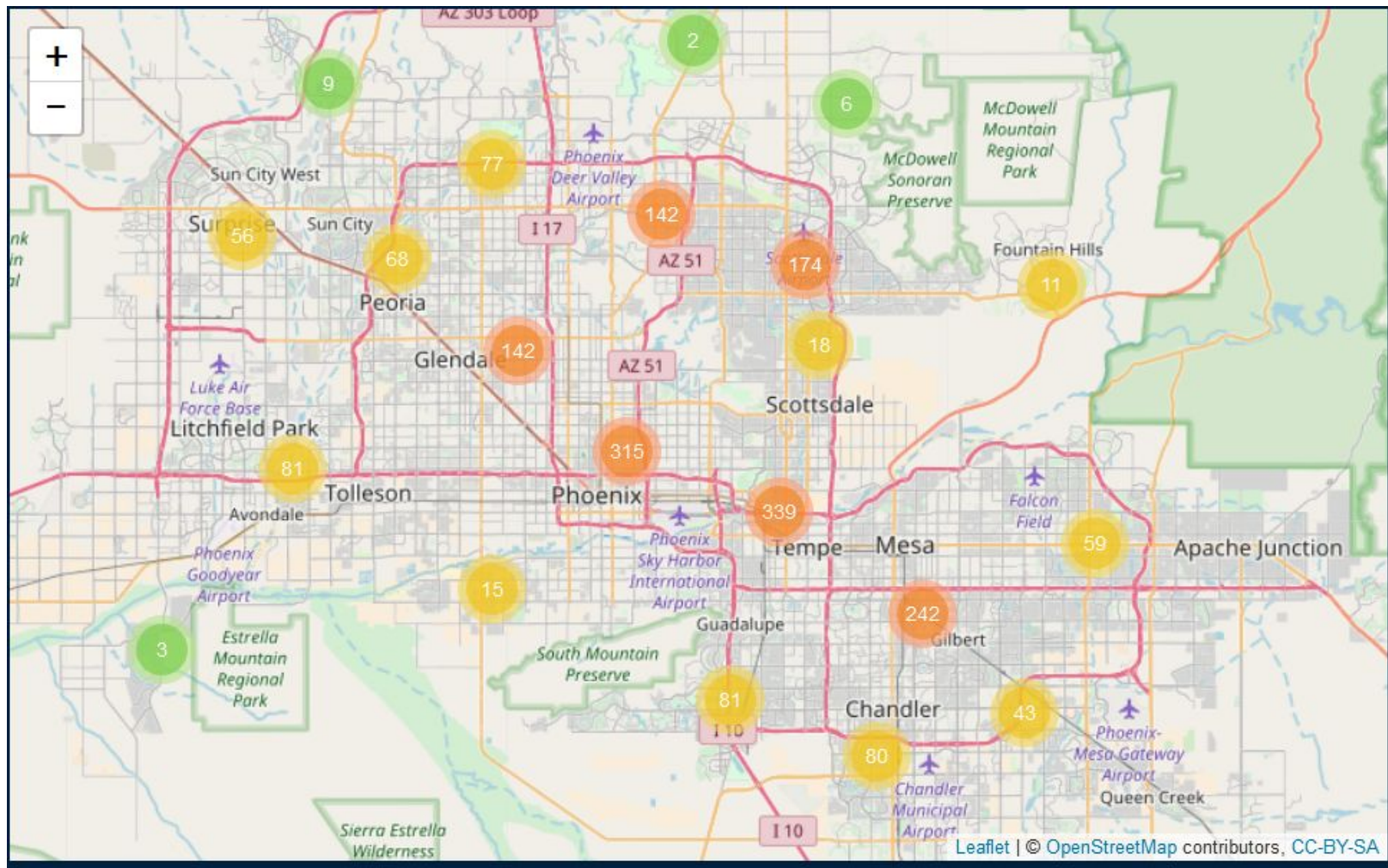
Descriptives



Descriptives





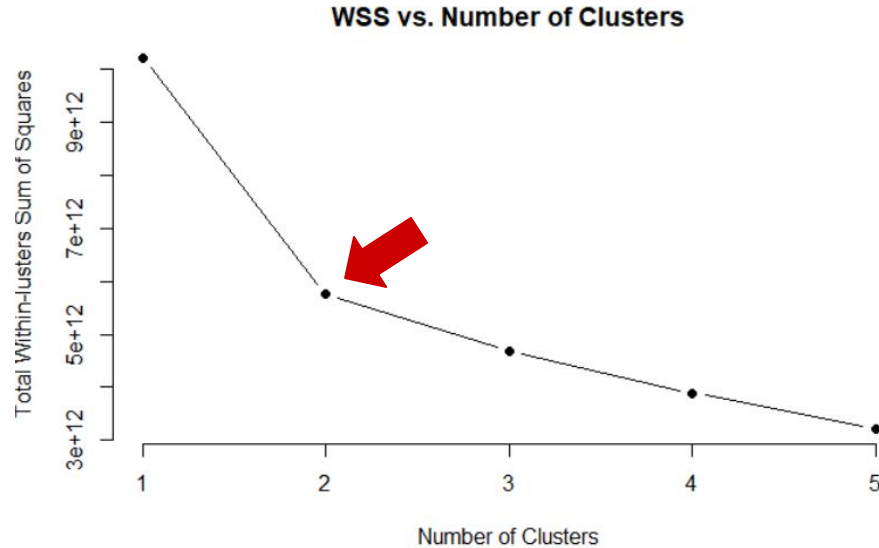


So now...the Machine Learning

- K-means Clustering
- Classification – two classes (≤ 3.5 stars vs. > 3.5 stars)
 - Logistic Regression
 - Quadratic Discriminant Analysis
 - Naive Bayes Classification
- Decision Trees

K-means Clustering

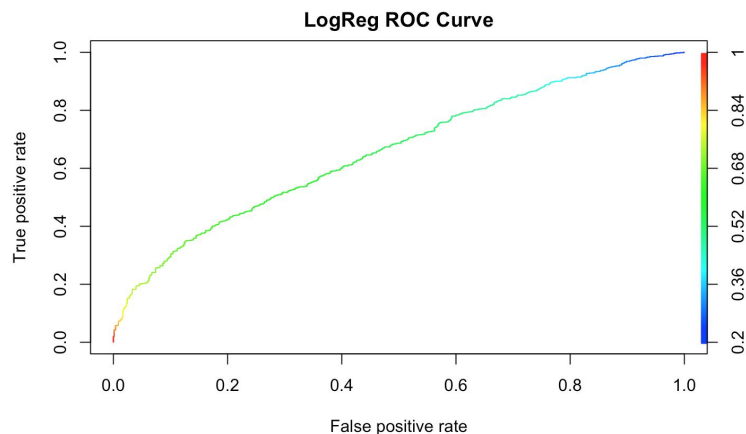
- Optimal number of clusters: $K = 2$
- Reinforces decision to use two classes for classification model



Logistic Regression & QDA

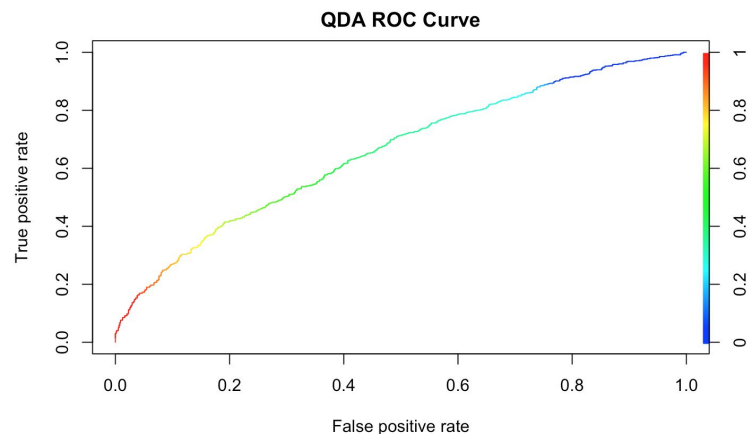
- Logistic Regression

- Overall lower test error
- Only correctly predicted **38.8%** of poorly-rated restaurants, but correctly predicted **79.1%** of well-rated restaurants



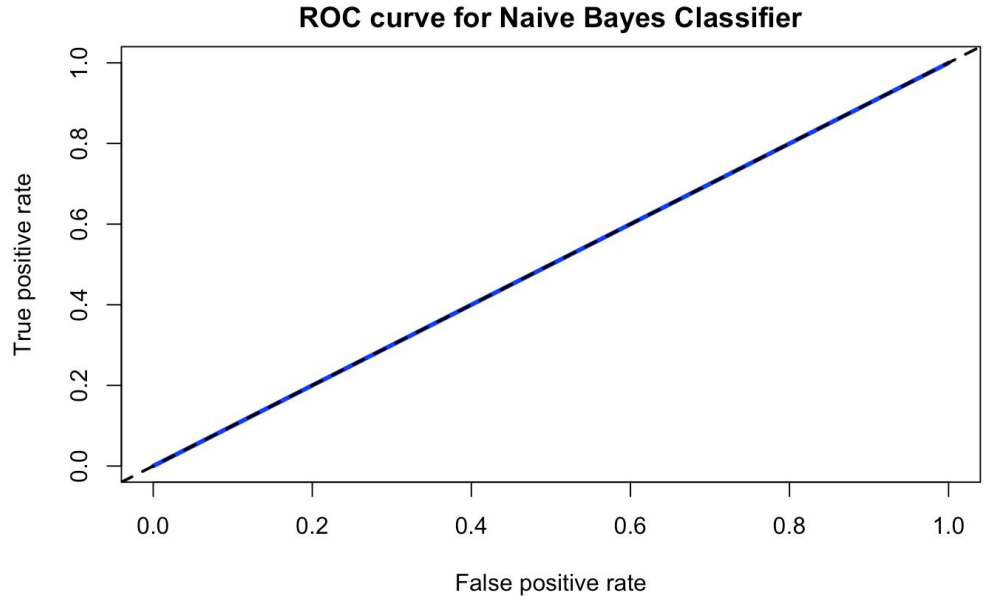
- QDA

- Slightly higher test error
- Better at predicting poorly-rated restaurants (**65.8%**), but not well-rated restaurants (**54.2%**)



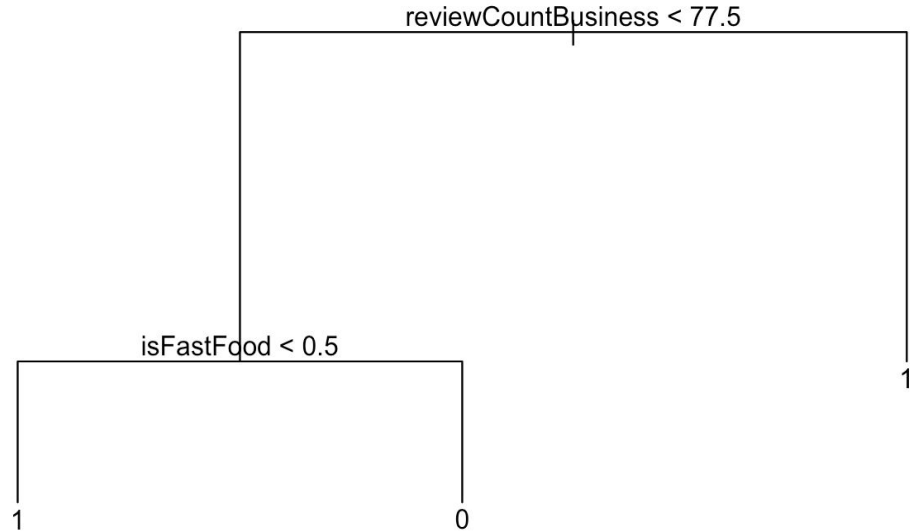
Naive Bayes for Classification

- Error rate is a little better than random guessing
- Incorrectly predicted poorly-rated
- AUC value: 0.5
- Picks majority class



Decision Trees

- Tree splits into only 3 nodes without the need for pruning
- Easy to predict rating based upon review count and whether or not the business is fast food
- Random forest model also clustered to majority class



Results

Technique	Test Error	Comments
Logistic Regression	0.3899371	<ul style="list-style-type: none">• Predicted highly-rated restaurants better
QDA	0.4059463	<ul style="list-style-type: none">• Predicted poorly-rated restaurants better
Naive Bayes	0.4894225	<ul style="list-style-type: none">• Predicted all as majority classifier

Conclusion & Future Work

- Every city is different
 - Would need to run models on every city to get most accurate results
- Naive Bayes did not perform well
 - Sentiment analysis could yield more insight into classification
- Important variables:
 - *reviewCountBusiness, avgDaysSinceJoined, isFastFood, isAsian, isAmerican, isBreakfast*
- More specific business attributes & demographics
- Define a better metric of success rather than star rating



Got questions?

