

# Raleigh Restaurant Inspections

JHUAPL

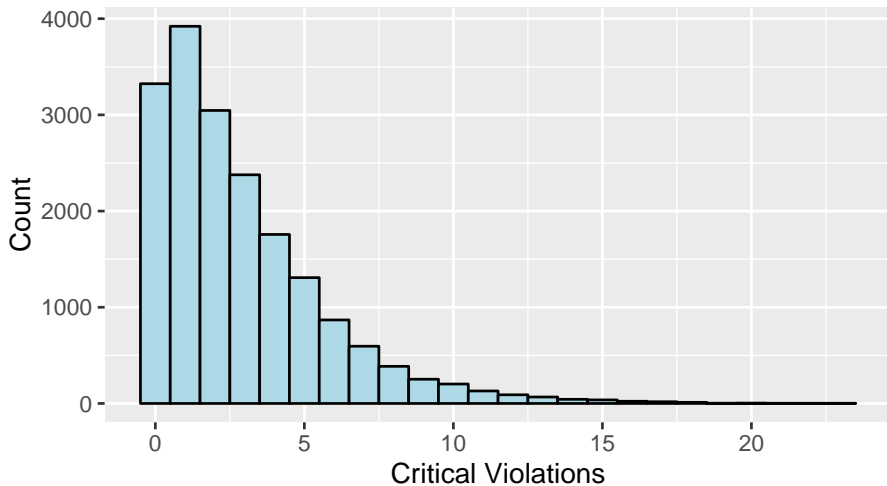
November 17, 2016

# Data Description

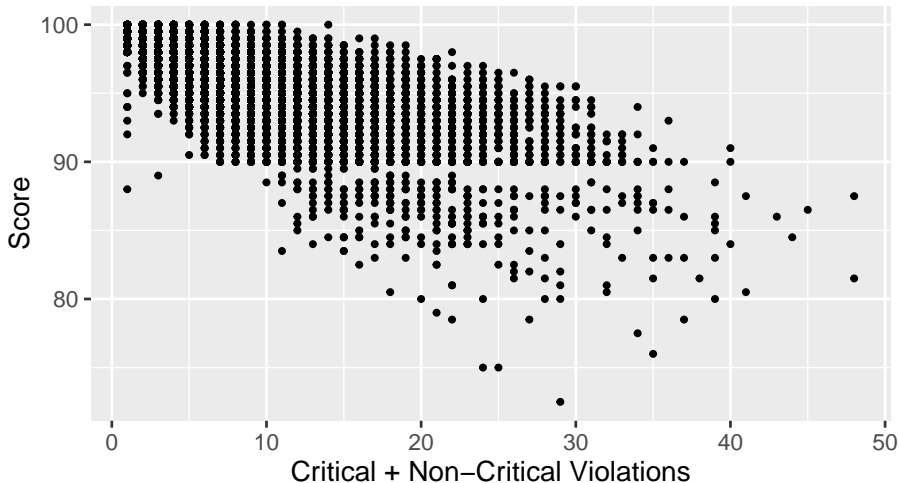
- Data from 9/21/2012 to 11/03/2016.
- 2,809 facilities (1,867 are restaurants)
- 18,469 inspections
- Cities in Wake County
  - ▶ top 5: Raleigh, Cary, Wake Forest, Apex, Morrisville

# Number of Critical Violations

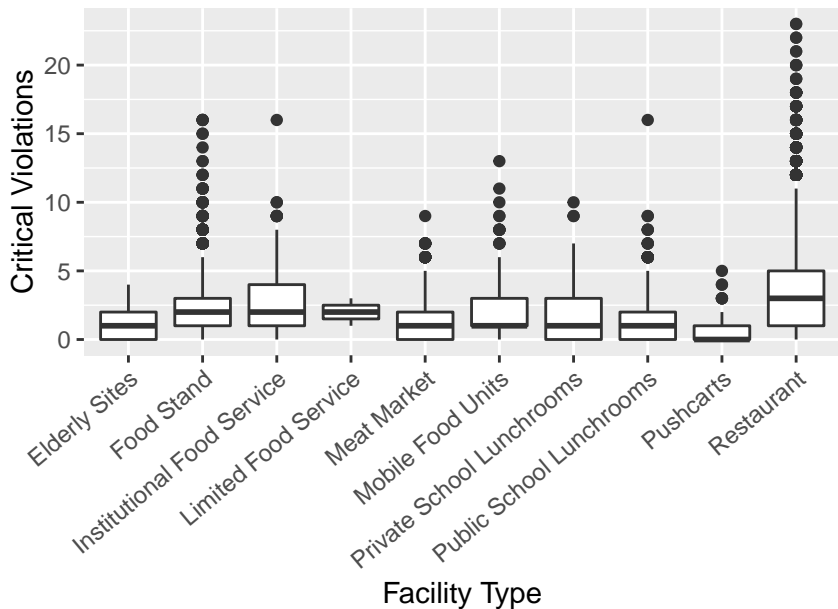
## Histogram of Number of Critical Violations



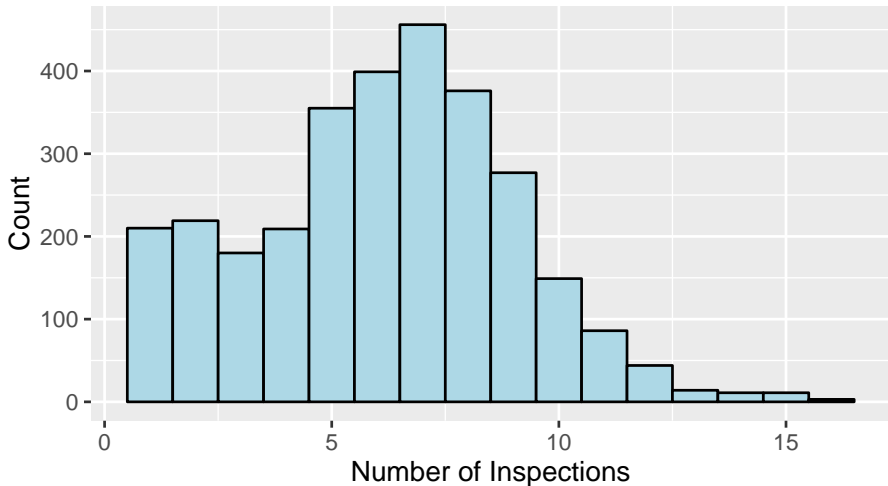
## Score vs. Number of All (Critical+Non-Critical) Violations



# Critical Violations by Facility Type



## Number of Inspections per Restaurant

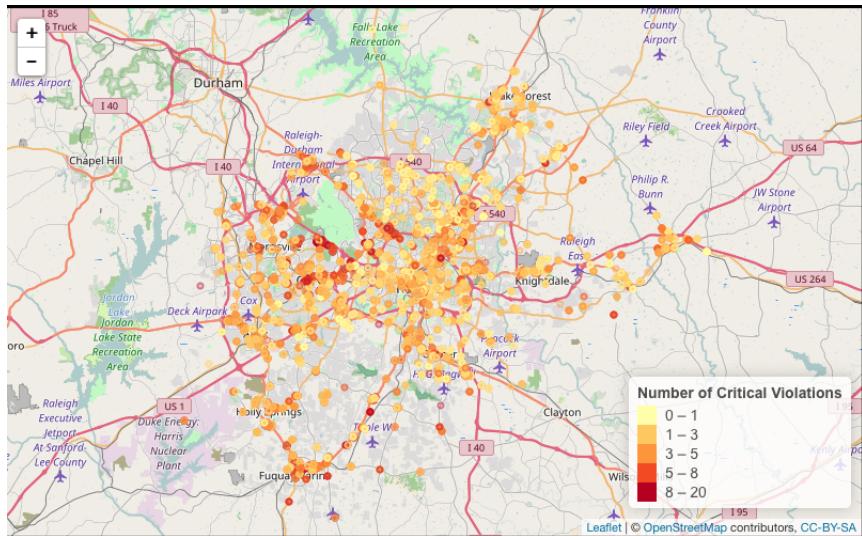


# Other Variables in Model / Other Data Sources

Besides number of previous critical violations, days from last inspection, and days since opening (extracted from inspections), other variables include:

- Geographical Information
  - ▶ Average Number of Critical Violations for all Prior Inspections of Nearest 5 Neighbors
- Census data: extracted income information by ZIP
  - ▶ Median Household Income
  - ▶ Percent below Poverty Line
- Yelp data: extracted information on restaurants
  - ▶ Rating (out of 5 stars)
  - ▶ Price (\$-\$\$\$\$)
  - ▶ Restaurant Category (top 20, e.g. Mexican, Sushi, Chinese, Coffee)

What is NOT in model: Inspector information, ZIP code

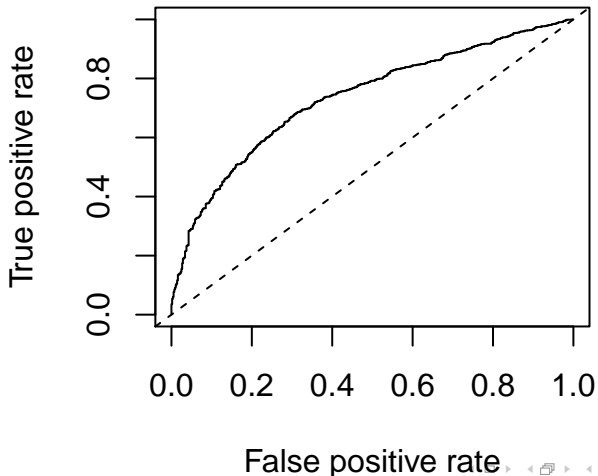






Logistic regression model trained on 2012-2015, tested on 2016 (Jan-Nov), approx. 70/30 split. The AUC on test-set is 0.734.

## ROC

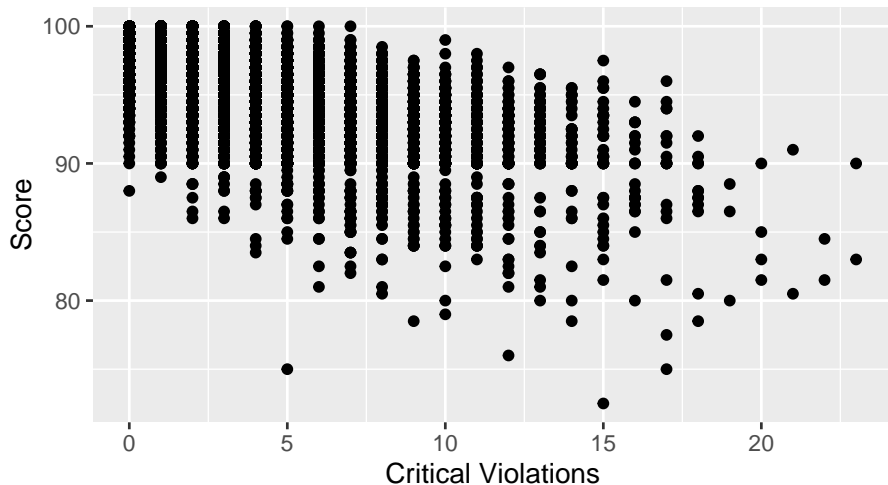


# Next Steps

- For binary response, what threshold would be most useful?
  - ▶ Results for cutoffs of 2, 3, 4, or 5 are similar (AUC 0.73-0.77) - beyond that gets rare
- Class imbalance
- To model counts, use a different model and performance measure
- How to make use of items that don't match from Yelp, Google Places
- Previous values seem quite useful: can we model those better (not just use one previous but all previous)

# Appendix A

## Score vs. Number of Critical Violations



# Appendix B

Variable Importance Plot (Random Forest)

avg\_neighbor\_num\_critical  
days\_from\_open\_date  
days\_since\_previous\_inspection  
num\_critical\_previous  
Median\_household\_income\_dollars  
Percent\_Families\_Below\_Poverty\_Line  
rating  
price  
grocery  
sandwiches  
hotdogs  
mexican  
pizza  
coffee  
tradamerican  
burgers  
chicken\_wings  
chinese  
delis  
breakfast\_brunch  
italian  
salad  
bbq  
seafood  
bakeries  
newamerican  
sushi  
bars

