



Group 5

TLC Taxi Dataset

A comprehensive view of vehicle for hire data in NYC

Mike Huffman, Tommy Guess, Joe Froelicher

The TLC Dataset

- Trip records from Yellow Cab, Green Cab and Ride Share Taxis for each month
- Attributes include pickup/dropoff location (by zone), time, trip distance, number of passengers, fare amount, tipping, etc.
- The dataset is split into month-specific CSV files. Depending on the area of study, data preparation varies.

Cabs of NYC



	Yellow Cabs	Green Cabs
Total Rides:	83,219,960	5,991,544
Manhattan	91.70%	34.7%
Queens	6.89%	29.5%
Brooklyn	1.21%	28.7%
Bronx	0.19%	7.0%
Staten Island	0.004%	0.05%
Newark	0.01%	0.002%

Questions to Answer

- How are rides distributed, geographically?
- How can we figure out the relationship or “association” between pickup regions and drop-off regions?
- What factors influence tipping?
- How did COVID impact taxi rides in NYC?

Summary Matrix, S

Algorithm 1 Building a Summary Matrix (S) from Rideshare Data

Require: X is an array of rideshare data where X_k is the row vector representing the k th ride and all of its associated statistics. DO is a feature of that set (and column of X) representing the drop-off zone, such that $DO_k \in \{1, 2, \dots, 265\}$ and represents the drop-off zone of the k th ride. PU is a feature of that set representing the pick-up zone, such that $PU_k \in \{1, 2, \dots, 265\}$ and represents the pick-up zone of the k th ride. S is initialized as the zero matrix $0_{265,265}$

for every ride X_k in X **do**

$DO_k \rightarrow i$

$PU_k \rightarrow j$

$S_{i,j} += 1$

end for

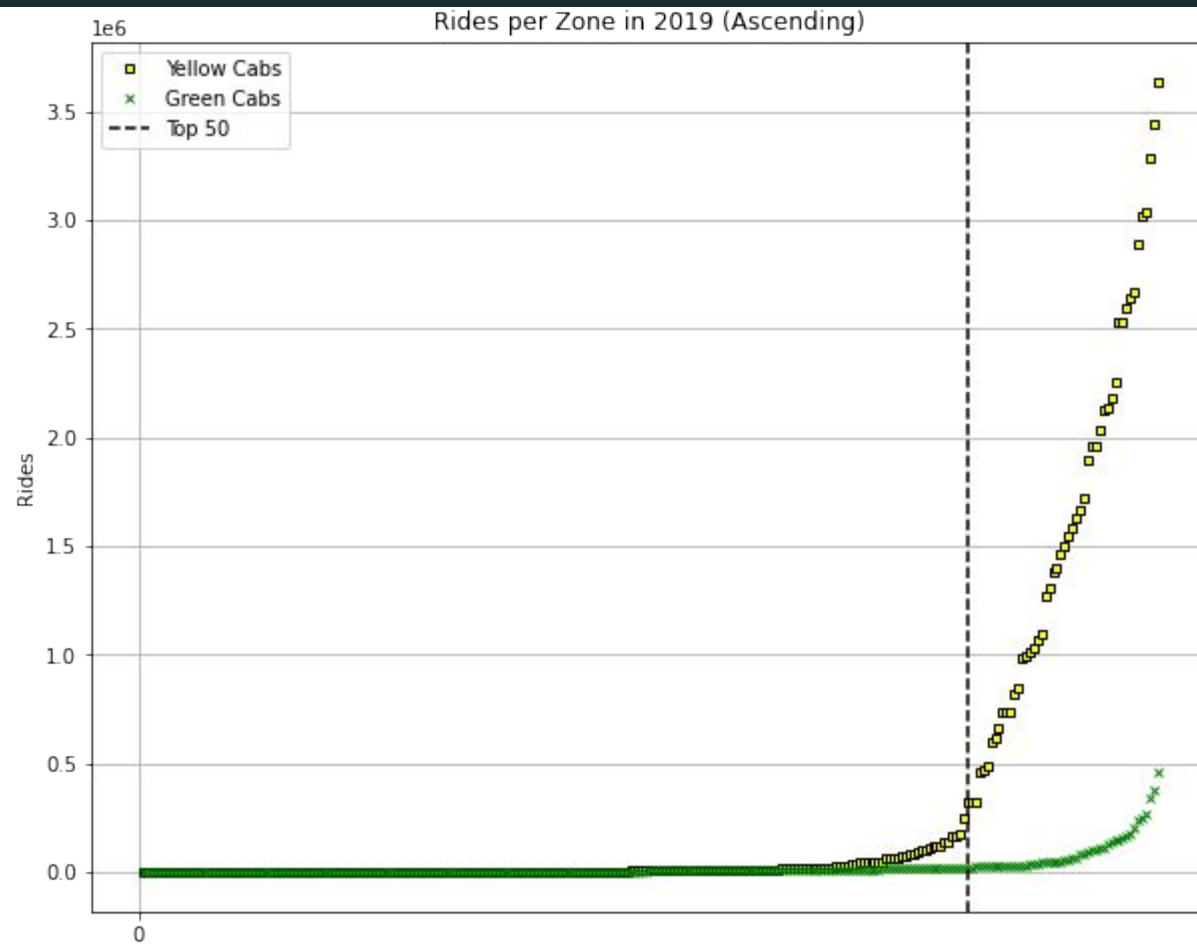
return S

Ride Coefficient (r.c.)

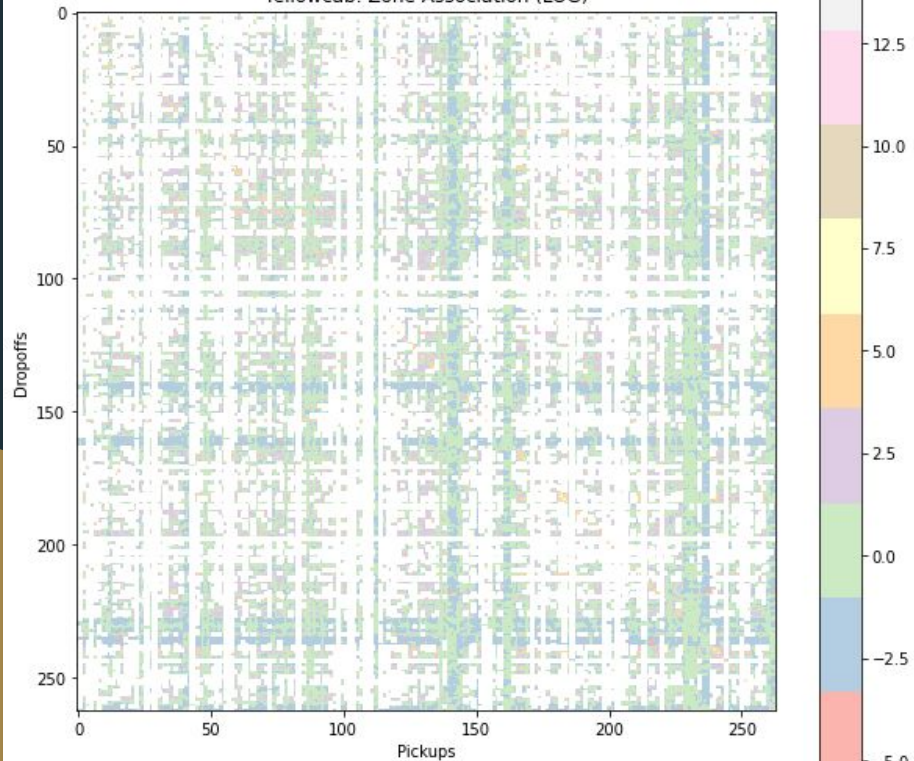
Let A be the event that ride pickups in zone x , and B be the event that the ride terminates in zone y .

$$r.c. = \frac{P(B|A)}{P(B)}$$

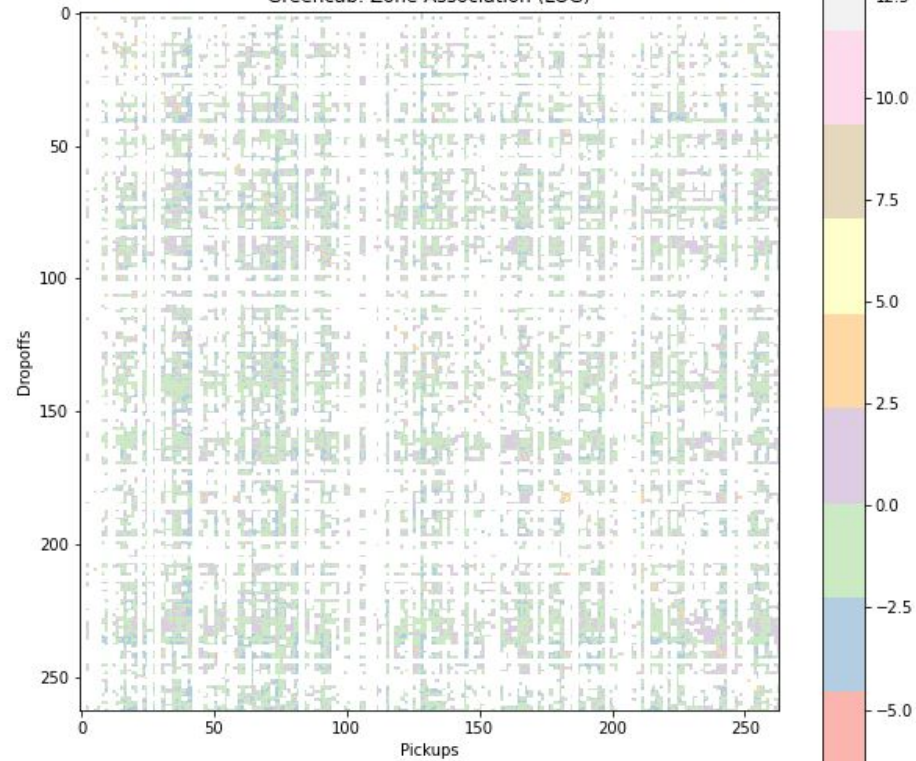
$$r.c. = \frac{P(A \cap B)}{P(A)P(B)}$$

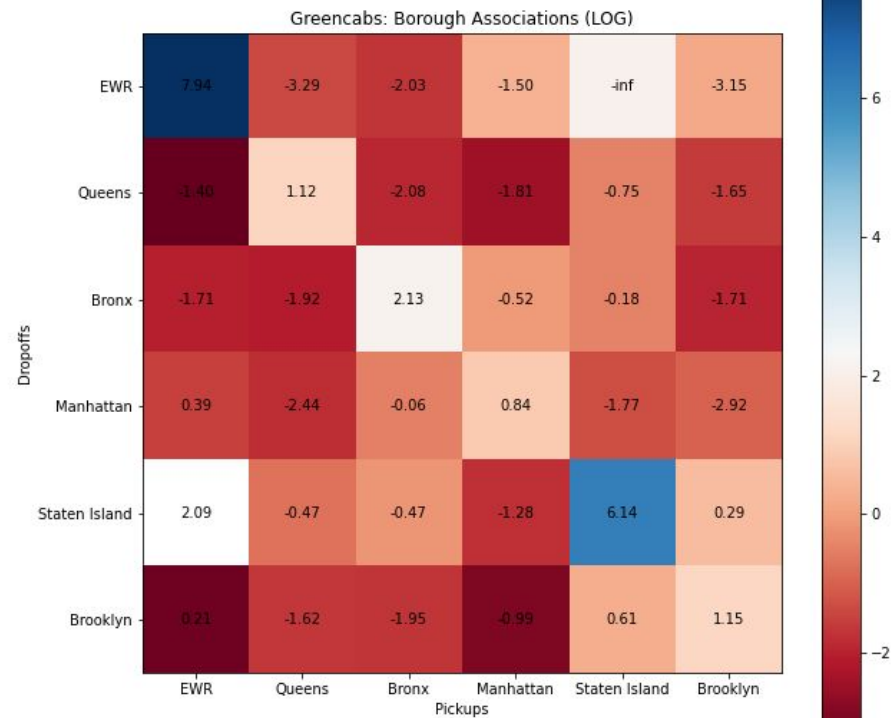
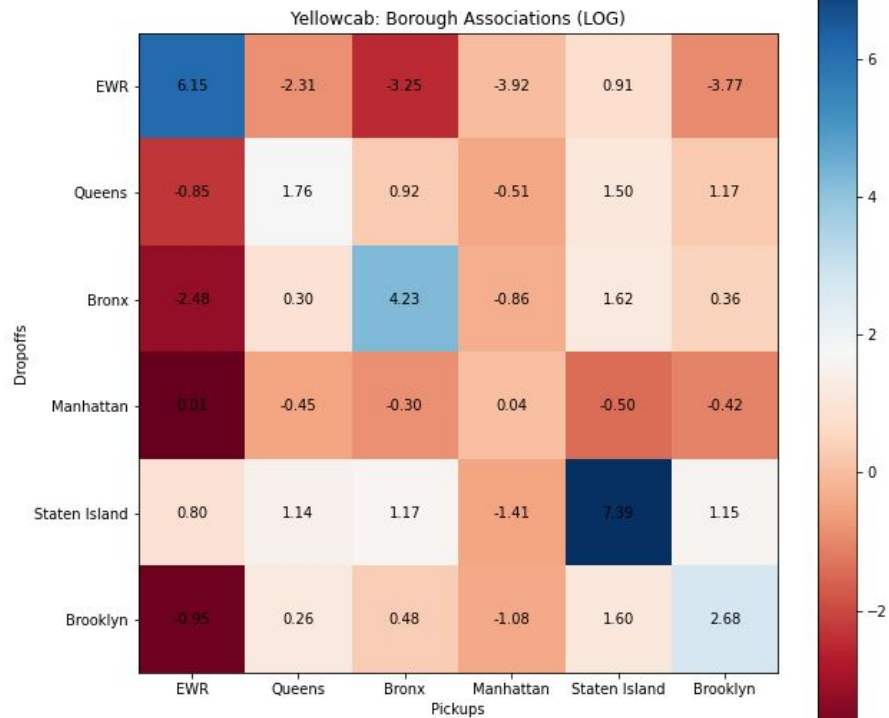


Yellowcab: Zone Association (LOG)



Greencab: Zone Association (LOG)



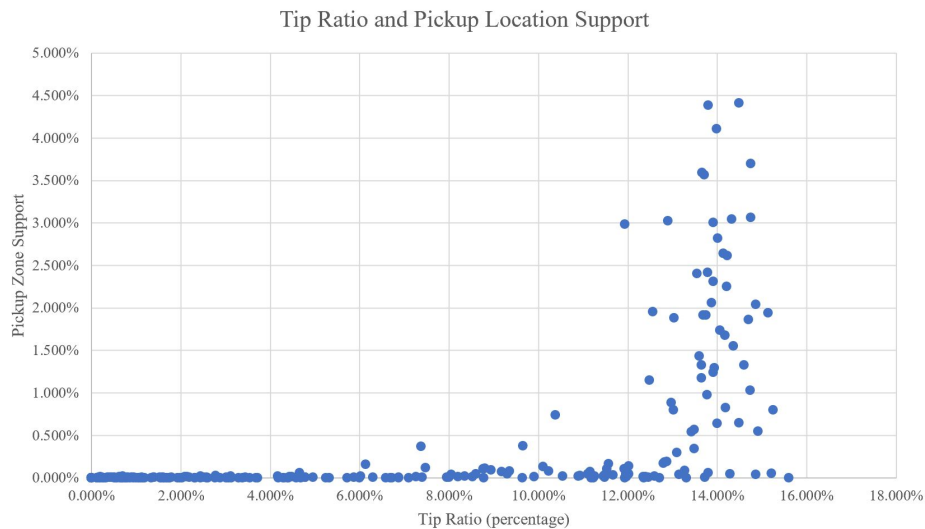


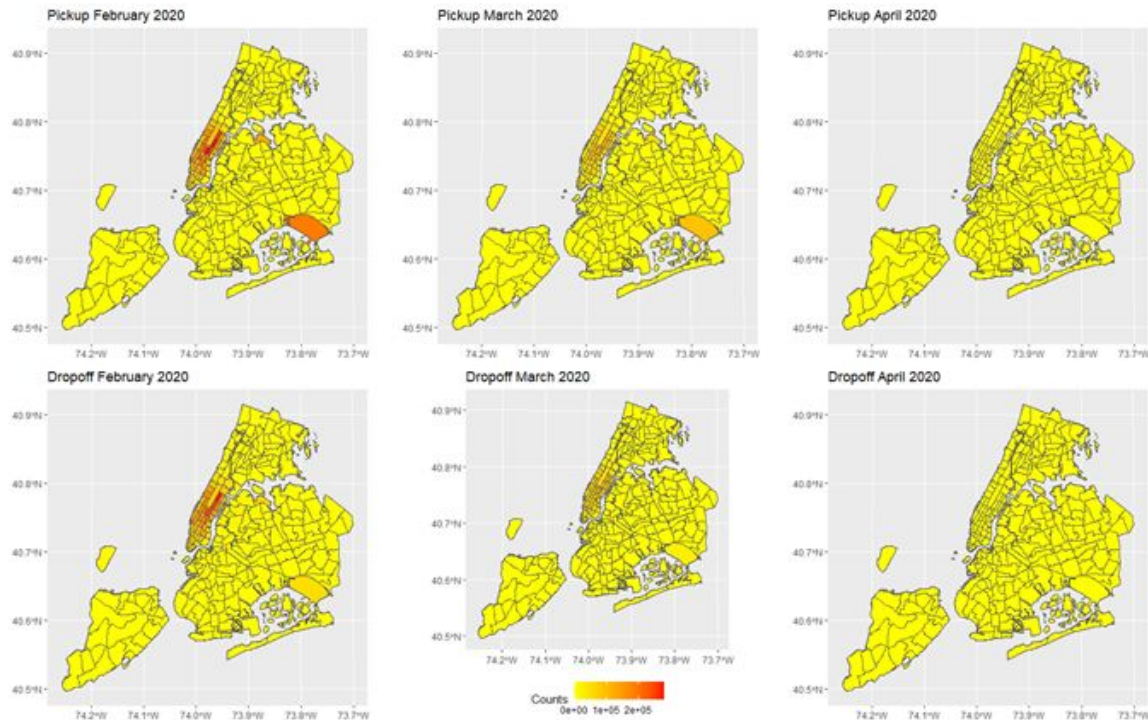
Financial Analysis

- Several attributes:
 - Tolls amount, improvement surcharge, MTA tax, etc.
- Most interesting: tips
- Rider has control over the tips
- Could potentially be a proxy for satisfaction
- Preliminary analysis included correlation matrix

Tip Ratio

- Tip as a percentage of fare





Maps: Normal cab patterns vs. abnormal patterns

What can we learn from the maps?

- Normal yellow cab traffic is primarily in Manhattan, and at La Guardia
- Color gradient from before the lockdown to during the lockdown
- Number of rides was negligible in those high trafficked districts

Covid-19 Effects on Yellow Taxi Rides

- Startling (but unsurprising) drop in the number of daily rides
- Normal weekly fluctuations diminished

What is the economic impact?

- Thousands of dollars change hands as a result of daily yellow cab rides

What other economic impacts might the fall of in yellow cab activity indicate?

- Where are people riding cabs going?
- Will they be spending money where they go?

