



The result of our permutation shows that there are no slopes that are steeper than our observed slope. All the gray lines of permuted regression slopes are less steep compared to our original slope.

```
In [32]: steeper_slopes = 0
n_permutations = 10000
permuted_slopes = np.zeros(n_permutations)

for i in range(n_permutations):
    fake_model = LinearRegression().fit(race_trip_station_df[['2017_median_income']], permute(race_trip_station_df[['2017_median_income']], permuted_slopes[i] - fake_model.coef_[0])
    permuted_slopes[i] = fake_model.coef_[0]

    if np.abs(fake_model.coef_[0]) > np.abs(model.coef_):
        steeper_slopes += 1

print("Percentage of slopes that are greater than our observed slope", steeper_slopes / n_permutations)
```

Percentage of slopes that are greater than our observed slope 0.0
Out of the 10000 iterations we created, we were never able to observe a permutation regression slope that was steeper than our data's original observed regression slope. This means that there is a near-zero probability of observing a slope this large under a null hypothesis. Therefore, we reject the null hypothesis. Thus, the conclusion we reached in Analysis Three shows a linear relationship between median income and trips per station, which is statistically significant and not due to random chance.

Interpretation and Conclusion

In attempting to answer our research question of whether Citi Bike is fairly distributing its stations in NYC, we performed three analyses.

First, to understand the influence of demand on station placement, we looked at the relationship between trips and stations. We observed a positive relationship between trips and stations in our scatterplot and more significant variability in trips as the stations increased. By plotting the location of geoids with the most stations on a map, we were able to find that geoids located near parks or docks had more stations regardless of the number of trips. In our linear regression, we numerically confirmed a moderately positive relationship between trips and stations with a Pearson correlation of 0.51 and a regression slope of 1.84e+04. Furthermore, we found that the R² score was 0.279, indicating that stations explained only 28% of the variability in trips. This led us to believe that other factors were influencing the number of stations in a census tract.

Thus, we investigated whether race and income play a role. In doing so, we found that census tracts with higher than the median number of stations were 77% White. On the other hand, census tracts below the median were 67% White indicating more stations were placed in Whiter areas. Moreover, through logistic regression, we found that the probability of a census tract having above the median number of stations increases with the proportion of White individuals. In contrast, this probability decreases as the proportion of Black and Asian individuals increases. Furthermore, through logistic regression between stations and income, we found that the likelihood of a census tract having above the median number of stations increases as the median income increases. Thus overall, we found that besides the number of trips, race and income also influence the number of stations in a census tract. More importantly, more stations are generally placed in White and higher-income census tracts.

To determine the fairness of Citi Bike's decision to place more stations in generally Whiter and higher-income census tracts, we looked at station utilization with race and income. Station utilization rate was defined as the ratio between the number of trips and stations in a census tract. We defined fairness as the placement of more stations in high-demand areas to ensure an equal utilization rate across all stations.

Comparing station utilization across all stations in majority Asian, Black, and White census tracts, we observed that on average, they had 30000, 6000, and 18000 trips per station, respectively. This relatively high station utilization rate in White areas justifies Citi Bike's decision to place more stations in Whiter census tracts, as they are just responding to demand. Moreover, the exceptionally high utilization rate in Asian census tracts, though it points to effective utilization, also implies a shortage of stations compared with the other census tracts. Predicting the future placement of Citi Bike stations, using linear regression, we observed that the number of stations in Asian census tracts would not increase significantly in the seven years from 2018 to 2025. Thus, if Citi Bike does not improve its current station placement strategy and begin focusing on Asian census tracts, this shortage will persist.

Finally, we analyzed station utilization across census tracts with different median income levels. We first found that census tracts whose median income falls below the 25th percentile have lower station utilization. Moreover, census tracts in the 75th percentile have a consistently high station utilization. Therefore, we suspected a positive relationship between income and station utilization and confirmed this through linear regression. The R² value of 0.17 suggests that 17 percent of the variability in the station utilization is explained by income. The Pearson Correlation, 0.41, and Spearman correlation, 0.51, indicates that the relationship between income and station utilization is moderate. Therefore, we conclude that Citi Bike is justified in placing more stations in higher-income census tracts as they have higher station utilization rates.

Overall, we have observed that Citi Bike's current placement of stations is not equitable, according to our definition of fairness. Additionally, our predictive analysis showed that Citi Bike would fail to address these issues in the years following 2017 if they remain complacent with their current strategy. Therefore we recommend Citi Bike address the problems we detailed above by increasing the availability of stations in Asian census tracts. If we had more time, we could have used other metrics to evaluate fairness, such as determining whether price is a prohibitive factor for Citi Bike utilization. Ideally, to do so, we would need a population survey detailing a person's income, ethnicity, and the maximum price they would be willing to pay for Citi Bike. Furthermore, if we could delve deeper into the topic, we would bring other data sets on road quality and crime rates in different census tracts. These hypothetical data sets would then be used to gain insight into how safety concerns affect a person's decision to utilize Citi Bike.

Limitations

Overall Limitations

- We acknowledge that race and income are not the only factors contributing to the distribution of Citi Bike stations around NYC.
- Since our data were based on Citi Bikes in 2017, our results are limited to that year only. That being said, they may be an indicator of how Citi Bike stations are distributed for other years.
- The data limits our definition of fairness in this project. There may be other definitions of fairness that we cannot evaluate in this project.
- Using census tracts as the geographical standard for dividing across different regions in NYC may have affected our results. If we used smaller geographical metrics (such as streets), we might have gotten more accurate and insightful results.

Citi Bike Trip Limitations

- The trip data only includes the Subscribers, which is about 90% of the original total trip data. Though there may be some one-time customers from NYC residents, due to a large amount of data, we decided to settle on subscribers since there is a higher chance that these individuals were New York City residents. So our data may not be representative of other NYC residents that may have only used one-time passes.
- We dropped rows where the start station equals the end station. When this occurs, we assume that most of these cases are people having trouble using Citi Bike initially. However, there may be cases where people made round trips and returned them to the same location for some of the trips where the start and end stations are the same.

Citi Bike Stations Limitations

- The station start date is not the date the station was added to the network but instead the date at which the first trip was taken from the station, which might not always be the same. However, this start date was used as a proxy for the date in which Citi Bike stations were made in our linear regression to predict future numbers of stations in each majority race census tract.

Race and Income Limitations

- There are only three racial categories and one additional 'other' racial category in our data. This indicates that some ethnicities are not reflected in the data.
- Populations like the homeless or digital nomads may not be recorded in the data.

Source Code

1. Team Github <https://github.com/shijessie/citibike-nyc>
2. Team Raw Data Google Drive https://drive.google.com/drive/folders/197_d1BwPHXmVuyKA6I7UEyEw4Nqo8dAK
3. Team Final Data Google Drive <https://drive.google.com/drive/folders/1gld5ily1ABeJ2GDxz-c-IUMjpu4Hktwa?usp=sharing>

Acknowledgements

- 1) MIT Citi Bike study : <https://aberke.github.io/income-race-bikes/>
- 2) Equity in Citi Bike study : <https://lrecr.pdx.edu/research/project/884>
- 3) Lecture Notes: https://colab.research.google.com/drive/1pMK18-DQoBDA_9rSzio6LJtkCJWxJF1g?usp=drive_open#scrollTo=6p2-tZKgm4He
- 4) Stack Overflow, Pandas, Matplotlib Documentation :
 - https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.xlabel.html
 - <https://stackoverflow.com/questions/11346283/rename-columns-in-pandas>

Special Thanks to Professor Wilkens, Head TA Stephen Cowpar, and all the TAs we have met in office hours and interacted through Ed in helping us answer our research question properly, find methods to solve our problems, and refine the direction of our project :)

Appendix

Citibike Trip DataCleaning Notebook : https://github.com/shijessie/citibike-nyc/blob/master/project_phase_submissions/Data_Cleaning.ipynb