

1 **ESTIMATING THE EFFECT OF UBER & LYFT ON PARKING VIOLATION IN NYC**

5 **Junjie Cai**

6 Institution: Center for Urban Science + Progress, New York University

7 Address: 370 Jay Street, Brooklyn, NY 11201

8 Email: jc9033@nyu.edu

10 **Junru Lu**

11 Institution: Center for Urban Science + Progress, New York University

12 Address: 370 Jay Street, Brooklyn, NY 11201

13 Email: lj1230@nyu.edu

15 **Yuxuan Wang**

16 Institution: Center for Urban Science + Progress, New York University

17 Address: 370 Jay Street, Brooklyn, NY 11201

18 Email: yw1665@nyu.edu

20 **Pranay Anchan**

21 Institution: Center for Urban Science + Progress, New York University

22 Address: 370 Jay Street, Brooklyn, NY 11201

23 Email: pya209@nyu.edu

25 **Shijia Gu**

26 Institution: Center for Urban Science + Progress, New York University

27 Address: 370 Jay Street, Brooklyn, NY 11201

28 Email: sg5718@nyu.edu

30 **Zhan Guo (Corresponding Author)**

31 Position: Associate Professor of Urban Planning and Transportation Policy

32 Institution: Wagner Graduate School of Public Service, New York University

33 Address: The Puck Building, 295 Lafayette Street, Room 3010, New York, NY 10012

34 Email: zg11@nyu.edu

37 Word Count:

38 words + 3 table(s) \times 250 = 750 words

45 Submission Date: August 2, 2019

1 ABSTRACT

2 This paper aims at exploring one potential ride-hailing services impact, taking the two biggest
3 ridesharing companies Uber & Lyft as examples: whether daily Uber & Lyft trips affect New
4 York City (NYC)'s parking demand, or specifically parking violations. NYC daily Uber & Lyft
5 trips, parking tickets data and some additional data were collected and aggregated by taxi zones.
6 Three technical models, Bayesian Network, Fixed Effects and Difference in Difference (DID) were
7 applied on these data. The results of these models showed a negative correlation and causal effect
8 between the number of Uber & Lyft trips and parking tickets, suggesting that Uber & Lyft may
9 help reduce parking violations and likely parking demand in NYC. Given the controversial issues
10 around ride-hailing services, this paper sheds light on the impact of Uber & Lyft and offer policy
11 insight to the NYC Taxi and Limousine Commission (TLC) regulation.

12

13 *Keywords:* Ride-hailing services, Parking Violations, Causal Inference, Bayesian Network, Fixed
14 Effect Model, Difference in Differences, New York City

1 INTRODUCTION

2 Ride-hailing services, like Uber & Lyft, are taking over marketing shares from Taxis (1).
3 They together have earned a revenue of \$13.46 billion in 2018 (2), while the revenue in U.S ride-
4 hailing segment is roughly \$44.8 billion (3). Nevertheless, the impact of Uber & Lyft is not limited
5 to the For-Hire vehicles (FHV) market alone. As a critical part of the transportation system, they
6 are reshaping our urban settings; although these effects are being questioned. On the one hand,
7 Uber & Lyft have created thousands of new jobs with an hourly salary between \$8.55 and \$11.77 on
8 average (4). On the other hand, they are suspected of increasing major city congestion significantly
9 (5), and even damage the public transit system of United States (6).

10 However, due to the limited availability of Uber & Lyft trip data (Uber & Lyft only have
11 public trip data available in NYC) and related scientific research, most of the Uber & Lyft effects
12 have not been uncovered and proven. To impose effective regulations on ride-hailing services, it is
13 critical for city administrators to understand them comprehensively. So, the goal of this paper is to
14 investigate one potential ride-hailing services impact – parking demand.

15 Theoretically, ride-hailing services may affect parking demand through several channels.
16 They may attract customers from private vehicles, traditional taxi cabs, transit systems, or people
17 who would forge their travel if these service were not available. The first scenario may reduce
18 parking demand as passengers do not need to park their cars. The second one may have little
19 effect on parking, while the 3rd and 4th scenarios may potentially result in more parking demand
20 as they are extra new vehicle travel even the parking need from rider-hailing service are relatively
21 minimum. The net outcome depends on the relative proportion of the four scenarios, so parking
22 demand might increase, decrease, or remain the same.

23 Note that the actual parking demand is hard to observe so a proxy measure is adopted
24 by this paper – parking violations. Parking violations are only applicable to street parking and
25 for a small portion of the parked cars, but it provides a proxy to understand parking demand.
26 This paper examines whether Uber & Lyft increase/reduce parking violations in New York City
27 (NYC). It leverages the power of large-scale machine learning and traditional statistics to model
28 and empirically prove the causal effect of Uber & Lyft rides on parking violations in NYC and
29 shed light on the impact of ride-hailing services on parking demand.

30 LITERATURE REVIEW

31 In a 2013 paper, Henao used a self-collected dataset containing information of 311 surveys
32 to ridesharing passengers. He asked a few questions in the survey, including the driving frequency
33 and the purpose of taking a ridesharing trip. By analyzing the dataset, he found out there are lots
34 of frequent drivers, who take TNC rides to avoid parking (7). This could be a potential factor in
35 the decrease in parking violations in the city. However, due to the size and nature of his dataset,
36 his analysis cannot statistically prove this effect.

37 Looking at patterns of parking violations in NYC, 'tickets for stopping or parking in illegal
38 zones, blocking traffic, tend to be most in common in areas with a greater commercial concentra-
39 tion of traffic.' The violations usually occur as a 'result of a driver needing to stop temporarily in
40 a busy area where parking is scarce; they try to park in a loading zone or double park for a few
41 minutes to avoid having to search for a parking spot.' (8)

42 As mentioned above, the current research mainly focuses on the study of ride-hailing ser-
43 vices or parking violations alone, but not to the extent of the relationship between them. Even if
44 some scholars is trying to fill the gap, they are frequently limited by the small size of the data.

This paper proposes a strategy combining large-scale machine learning methods and traditional statistics with several open datasets to effectively infer the causal effect on the parking violations brought by ride-hailing services in NYC.

Basically, causal inference is the process of drawing conclusions about a causal connection based on the conditions of the occurrence of an effect. In Pearl's 2009 paper, he discusses the advances in causal inference and "the approaches to be undertaken when moving from traditional statistical analysis to causal analysis of multivariate data." He surveys the "development of mathematical tools for inferring answers to three types of causal queries: queries about causal effects, probabilities of counterfactuals and direct and indirect effects (mediation)." (9) He discusses how "combining the algebraic component of the structural language of causality along with the graphical component provides a powerful and comprehensive methodology for empirical research" (9). In this paper, specifically three models - Bayesian Network, Fixed Effects and Difference-in-Differences (DID) are adopted.

In a 1996 study, it was shown how the "use of causal independence in a Bayesian network could greatly simplify probability assessment as well as probabilistic inference." (10) Heckerman and Breese look at an important weakness in a real-world Bayesian network representation, i.e. an effect could have potentially many causes; probability specification and inference could be impractical or impossible. They find some "potential sources for gains in inference efficiency, although these gains are not very apparent for more general Bayesian networks." They used a few models on several artificial and real-world Bayesian networks to better understand general Bayesian networks. In each model, they used the noisy-MAX model (11) to encode all parent-child relationships. They "transform Bayesian networks into an annotated undirected tree where each clique (node) corresponds to a set of nodes in the original network". Their results indicate the "use of decomposition can decrease inference complexity substantially when nodes have many states and many parents."

Fixed effects models provide a way to estimate causal effects in analyses where units are measured repeatedly over time. It can eliminate the effects of confounding variables without measuring them or even knowing exactly what they are, if they remain constant over time. One significant drawback with this model is that 'a great deal of information can be lost by focusing only on variation within individuals, thereby ignoring the variation across individuals.' Since this model removes the effects of all time-invariant causes, the standard fixed effects model is unable to estimate the effects of time-invariant measured causes (12).

Difference-in-differences (DID) is another mainstream statistical technique widely applied in econometrics and social sciences for causal inference. DID attempts to use observational data to mimic the design of experimental research, studying the differential effect of a treatment. (13) Alberto Abadie from Harvard University is one of the representatives who applied the DID model to the field of sociology. He used this idea in 2010 to effectively estimate the negative impact of California's Proposition 99 on local tobacco sales (14).

DATA

This paper uses three datasets to quantify the causal effect of estimating ride-hailing services versus parking tickets. All of them are aggregated with taxi zones, which is an officially defined geographical unit that roughly based on NYC Department of City Planning's Neighborhood Tabulation Areas (NTAs) and are meant to approximate neighborhoods. See the general attributes of taxi zones at <https://data.cityofnewyork.us/Transportation/NYC-Taxi-Zones/d3c5-ddgc>.

1 **Parking Ticket Data**

2 The parking violations data were obtained from NYC Open Data, and this project used data
3 from 2014 to 2018. NYC issues about 11 million parking tickets each year. Four columns, Issue
4 Date, Violation Code, Street Name, and Borough, and 80 types of parking tickets, which related
5 to private vehicle, are used. This parking tickets data is completely open and being updated daily,
6 which can be referred at this website: [https://data.cityofnewyork.us/City-Government/
7 Parking-Violations-Issued-Fiscal-Year-2014/jt7v-77mi](https://data.cityofnewyork.us/City-Government/Parking-Violations-Issued-Fiscal-Year-2014/jt7v-77mi).

8 Since the dataset only provides borough and street name for each parking ticket, Google
9 API was used to geocode the given borough and street name for each observation to an approximate
10 location with longitude and latitude and then mapped on to according taxi zones. Eventually, the
11 dataset was grouped by date and taxi zone to show the number of parking tickets issued in each
12 taxi zone in New York for each day from 2014 to 2018.

13 **Uber & Lyft Ride Data**

14 The Uber & Lyft trip data was extracted from the NYC For-Hire Vehicle (FHV) dataset
15 that is available on the Taxi & Limousine Commission (TLC) website: [https://www1.nyc.gov/
16 site/tlc/about/tlc-trip-record-data.page](https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page). Dispatching base number, which is provided
17 for each observation, is different for different FHV companies. So it was used to identify which
18 trips were conducted by Uber & Lyft.

19 The collected dataset for Uber & Lyft trips includes fuzzy location by taxi zone (there are
20 263 taxi zones in New York City, and the mean partition granularity is 3 square kilometer) and time
21 of pick up and drop off for each trip. Spark was used to group the dataset by trip date and pick-up
22 and drop-off taxi zones. Pick-up trips were finally used as the number of Uber & Lyft trips. Each
23 row of the prepared dataset indicates the number of trips conducted by Uber & Lyft originated
24 within one taxi zone for each day in the four-year window. Parking ticket data and Uber & Lyft
25 trips data were merged according to the taxi zone and the date to form a final prepared dataset.

Table 1 shows the prepared data.

TABLE 1 A sample of the prepared data

Taxi zone id	Date	Pickup counts	Ticket counts
1	2015-01-01	1	10
3	2015-01-01	9	0
4	2015-01-01	411	43
6	2015-01-01	2	22

27 **Other Data**

28 A general idea of causal inference research is to control third-party factors that uncorrelated
29 with the independent variable but correlated with the dependent one. Thus, five additional spatial
30 datasets that may affect the parking tickets from other aspects (geographical, demographic and
31 socio-economic) were included. They were used in the DID as clustering variables and in the
32 Bayesian Network as potential causation for parking violation (**Table 2**):

- **ACS data** was extracted from the American Community Survey 2015 5-year estimates.

18 statistical attributes, including population density, poverty rate, etc., were used to add the socioeconomic information for all models.

- **NYC crime data** is accessible on NYC Open Data. Three variables: number of felonies, number of violations, and number of misdemeanors, were mapped to taxi zone level and used in the model.
- **SAT result data** is hosted on NYC Open Data. The average score of SAT reading, math, and writing sections was included to measure the education level of each zone.
- **Transportation accessibility data** is available on NYC Open Data. This dataset includes subway entrance and bus stop coordinates. The number of subway entrances and the number of bus stops represent the transportation accessibility.
- **Parking capacity data** was collected from NYC Open Data, including meter parking and parking lot. They indicate the parking capacity of taxi zones.

TABLE 2 Composition of 5 additional datasets

Dataset	Attributes
ACS	'DensityPop': Population Density 'TotalPop': Total Population 'IncomePerCap': Income per capita (\$) 'Poverty': under poverty level rate (%) 'Professional': % employed in management, business, science, and arts 'Service': % employed in service jobs 'Office': % employed in sales and office jobs 'Construction': % employed in natural resources, construction, and maintenance 'Production': % employed in production, transportation, and material movement 'Employed': employed rate (16+) (%) 'Unemployment': Unemployment rate (%) 'Drive': % commuting alone in a car, van, or truck 'Carpool': % carpooling in a car, van, or truck 'Transit': % commuting on public transportation 'Walk': % walking to work 'OtherTransp': % commuting via other means 'WorkAtHome': % working at home 'MeanCommute': commute time (minutes)
Crime	'FELONY': Number of felony crimes in the taxi zone 'VIOLATION': Number of violation crimes in the taxi zone 'MISDEMEANOR': Number of misdemeanor crimes in the taxi zone
SAT	'sat': Average score of SAT reading, math, and writing
Transportation	'subway': Number of subway entrances 'bus': Number of bus stops
Parking	'meter': Number of meter parking 'parkinglot': Area of parking lot

1 METHODOLOGY

2 Three models - Bayesian Network, Fixed effects, and Difference in Differences (DID) were
 3 developed. The combination of the three methods was expected to comprehensively explain the
 4 causal effect of Uber & Lyft on parking violations. All of our data preparation and coding work can
 5 be reached at: <https://github.com/uberlyftparkingviolation/NYU-CUSP-Capstone-2019>.

6 Bayesian Network

7 The key idea of Bayesian Network is that one can distinguish correlation from causation if
 8 independent causes can be observed (9). This project integrated 14 observational features of taxi
 9 zones: 'Pickup', 'Ticket', 'PopDensity', 'IncomePerCap', 'Poverty', 'Professional', 'Employed',
 10 'Crime', 'Subway', 'Bus', 'Meter', 'ParkingLot', 'Sat', and 'Carpool'. Bayesian Networks can
 11 be learned from the datasets to demonstrate whether Uber & Lyft directly/indirectly cause parking
 12 violations in NYC. There are three main steps: structure learning, parameter learning, and link
 13 strengths measuring.

14 Structure learning estimates a directed acyclic graph (DAG) that captures the dependencies
 15 between the variables given a set of data records. This paper combines two algorithms to learn the
 16 DAG: score-based structure learning and constraint-based structure learning.

17 Score-based Structure Learning approach construes model selection as an optimization
 18 task. It has two building blocks: first apply a 'scoring function' $s_D: M \rightarrow \mathbb{R}$ that maps models
 19 to a numerical score, based on how well they fit to a given data set D ; second perform 'search strat-
 20 egy' to traverse the search space of possible models M and select a model with optimal score. For
 21 the first building block scoring function 'Bayesian Dirichlet', scores such as BDeu, K2, and BIC
 22 (Bayesian Information Criterion) were used to measure the fit between model and data. For the
 23 second building block 'Search strategy', a heuristic search 'HillClimbSearch' was used, which im-
 24 plements a greedy local search that starts from the DAG^{start} (disconnected DAG) and proceeds
 25 by iteratively performing single-edge manipulations that maximally increase the score.

26 Constraint-based Structure Learning attempts to correctly capture the directionality of causal
 27 relationships. It has two building blocks: identifying independencies in the data set using hypoth-
 28 esis tests and constructing a DAG according to identified independencies. For the first step, inde-
 29 pendencies in the data can be identified using χ^2 conditional independence tests. The \hat{p} -value is
 30 the probability of observing the computed χ^2 statistic or a higher χ^2 value, given the null hypothe-
 31 sis that X and Y are independent for given Z s. This can be used to make independence judgments,
 32 at a given level of significance: p -values less than the chosen significance level α reject the null
 33 hypothesis, and it can be concluded that X and Y are conditionally dependent given Z s. P -values
 34 greater than or equal to α fail to reject the null hypothesis, from which it can be concluded that X
 35 and Y are conditionally independent.

36 Parameter learning is a reinforcing stage of Bayesian Network method. The task of param-
 37 eter learning is to estimate the values of the conditional probability distributions (CPDs). Given a
 38 Bayesian network structure and a training dataset, one can learn the parameters (conditional distri-
 39 bution of each node for each distinct combination of parent values) by maximum likelihood. There
 40 are two main algorithms: Maximum Likelihood Estimation (MLE) and Bayesian Parameter Esti-
 41 mation. MLE simply uses the relative frequencies with which the variable states have occurred.
 42 The Bayesian Parameter Estimator starts with already existing prior CPDs, that express our beliefs
 43 about the variables before the data was observed. Those "priors" are then updated, using the state
 44 counts from the observed data. Parameter learning would be performed in the future to further

improve the causal inference ability of our network.

Pearson correlation coefficients are calculated to measure the strength of connection along a specific edge. Link strengths among socioeconomic features are represented by spatial correlations since they are differential though zones but rarely vary through years. Link strength between Uber & Lyft and parking violations is represented by the combination of spatial correlation and temporal correlation. For the time series correlation, we are focusing on the periodic variance rather than the constant part. Thus Hodrick-Prescott (HP) Filter is implemented to split the overall trend and cyclical pattern. The overall trends represent the fixed parts of tickets and Uber & Lyft trips that must come up no matter what happens. And the cyclical patterns represent the daily variations. Finally the temporal correlation was calculated with the filtered cyclical variations.

Fixed Effects Model

Fixed effects model (FEM) is widely used to control for unobserved variables in causal inference (15). FEM fits this study well, which aims to explore the causality using multidimensional datasets shown in **Table 1**. The datasets were standardized before modeling in order to improve the optimization training process and measure variable importance. To develop a FEM, dummy variables of each taxi zone were added to a standard OLS model to evaluate the fixed effects of the number of Uber & Lyft trips across different taxi zones. Thus the effect of Uber & Lyft trips can be separately measured after filtering the unobserved fixed effects. Considering the model for taxi zones $i = 1, \dots, N$ which is observed at certain time periods $t = 1, \dots, T$ (**Equation 1**):

$$y_{it} = \alpha + \beta \sum_{i=1}^N \sum_{t=1}^T x_{it} + \gamma \sum_{i=1}^N z_i + u_{it} \quad (1)$$

where y_{it} and x_{it} are, respectively, the number of parking tickets (dependent variable) and the number of Uber & Lyft trips (independent variables) in the i -th taxi zone at time t ; z_i are unobserved, time-invariant dummy variables in different taxi zones; β and γ represent coefficients of the independent variables and dummy variables; α and u_{it} refers to the intercept and the error term. The coefficient, β , of the independent variables would indicate the causal effect exerted by Uber & Lyft trips on parking violations.

DID: Difference in Differences

Within the framework of natural science, two identical NYCs, one having TNC companies operating in it and the other not, were needed to strictly prove the causal effect of Uber & Lyft on parking violations. However, such ideal 'treatment groups' and 'control groups' were impossible to be found or developed in social science research. Consequently, DID model was applied in this project to solve this problem. **Figure 1** shows the general idea of DID and how it works. The 'treatment object' and 'control object', P and S , have status P_1 and S_1 on time period 1, and status P_2 and S_2 on time period 2. Based on the DID's most critical assumption, parallel trend assumption, that any existing factor would have same effect on P and S , shown by the parallel S_1 to S_2 line and P_1 to Q dotted line DID, the difference between Q and S_2 on the second time period equals the difference between P_1 and S_1 on the first time period, and the difference between P_2 and Q represents the effect of the outside factor or treatment.

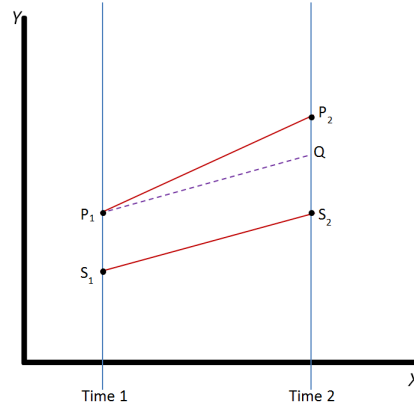


FIGURE 1 Illustration of DID model, by Danni Ruthvan

Average causal effect can be expressed by **Equation 2**:

$$\frac{1}{n} \sum_{i=1}^n \frac{\Delta Outside}{(P_2 - S_2) - (P_1 - S_1)} \quad (2)$$

where $\Delta Outside$ refers to the difference between the average amount of Uber & Lyft trips in 2015-2018 between the 'treatment taxi zone' P , and the 'control taxi zone' S , for every i . P_1, P_2, S_1 , and S_2 refer to the average amount of parking tickets in 2014 and 2015-2018 between P and S . n refers to the number of 'treatment vs. control' pairs.

To meet the parallel assumption, cross-clustering was applied to find multi-dimensionally homogeneous taxi zones to be used as 'treatment group' and 'control group.' The taxis zones were clustered based on two dimensions separately, extracting those zones always belonging to the same clusters: the number of parking tickets for each taxi zone in 2014, which, shown by data, had not been affected by Uber & Lyft; and the above 5 additional datasets, which summarizing different aspects of taxi zones from geographic, demographic and socio-economic views that will also hardly affected by Uber & Lyft over time. In each cluster, the remaining zones were pairwise paired, which generated over 900 pairs of homogenous taxi zones.

1 RESULTS

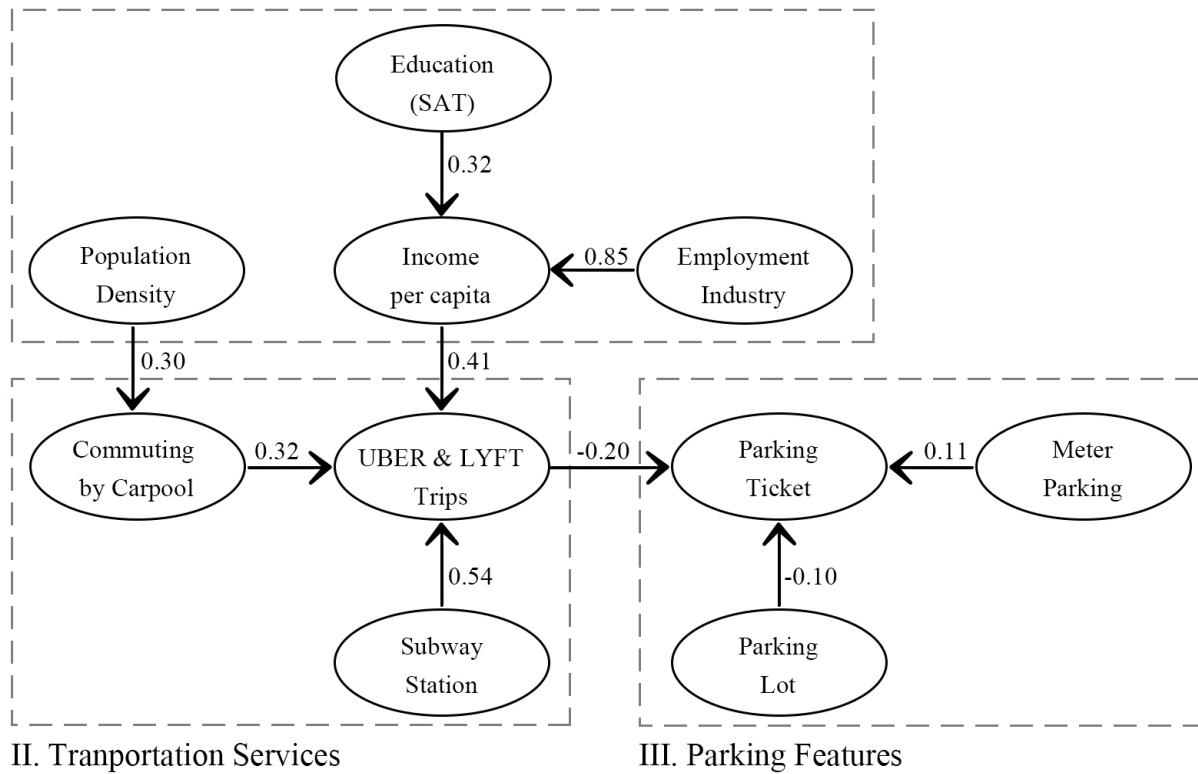
2 Results of these three models are listed below, which proves the negative casual effect on
3 parking tickets brought by Uber & Lyft trips in NYC.

4 Bayesian Network

5 The causal network structure generated by Bayesian Network indicates that Uber & Lyft
6 have a direct impact on parking tickets. According to **Figure 2**, there are three variables that
7 directly influence the parking ticket, and Uber & Lyft is the strongest one of them. After applying
8 HP filter (**Figure 3**), the Pearson's correlation coefficient is -0.20 between the cyclical variations of
9 Uber & Lyft trips and parking tickets. And p-value is $6.52 * e^{-14}$, much lower than the significant
10 level of 0.05. In conclusion, Uber & Lyft statistically significantly reduce the parking violations.

11 However, there are limitations and potential biases with Bayesian Network. Initially, Bayesian
12 Network method bases on many hypotheses, including causal Markov, causal faithfulness, causal
13 sufficiency, and acyclicity, but it is difficult to meet all of them in real world. Meanwhile, unob-
14 served intermediate factors may be missing between nodes. That is, Bayesian Network could never
15 be absolutely certain whether the influence is direct.

I. Demographic Properties

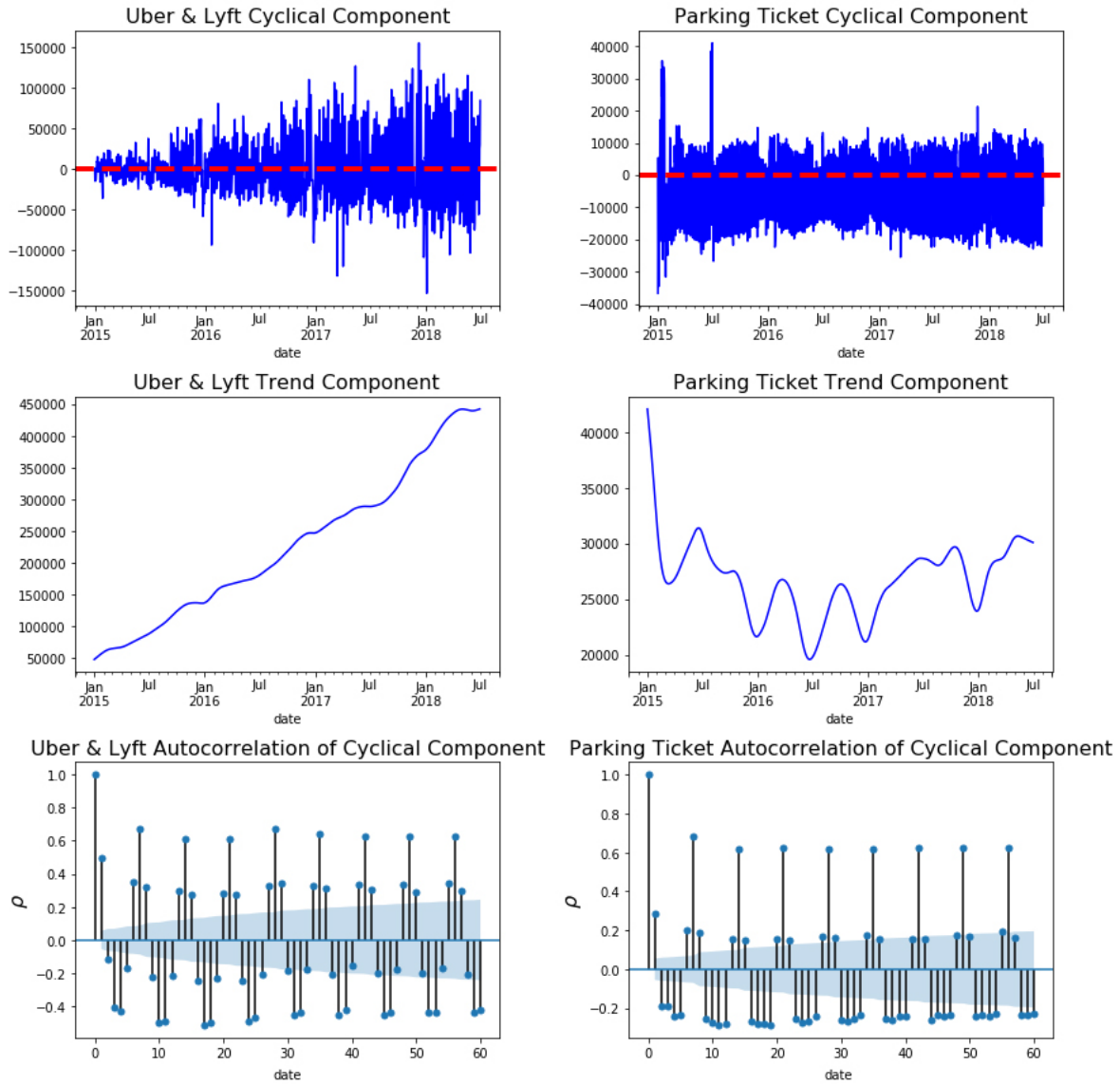


Notes:

Causal Network structure is learnt by Bayesian Network algorithm.

Association strengths among edges are represented by Pearson's Correlation Coefficients.

FIGURE 2 Bayesian Network

**FIGURE 3 Fixed Trend and Cyclical Variations Filter**

1 Fixed Effects Model

2 The fixed effects model (FEM) controls for the unobserved factors. According to **Figure**
 3 **4**, R-squared of the Uber & Lyft regression is low since the number of parking tickets cannot be
 4 explained by Uber & Lyft trips solely. But after adding dummy variables, the model achieves an
 5 R-squared of 0.70, well explaining the variation of number of parking tickets. The coefficient of
 6 Uber & Lyft trips (pickup) is -0.0007, indicating Uber & Lyft trips slightly reduce NYC parking
 7 violations. However, this effect is not statistically significant, because the p-value of is 0.677,
 8 much larger than the significance level of 0.05. And the confidence interval across zero (-0.0039,
 0.0026) shows the effect of Uber & Lyft is volatile through taxi zones.

Fixed Effect Model Result Summary						
=====						
Dep. Variable:	tickets	R-squared:	5.197e-07			
Estimator:	Pooled Least Squares	R-squared (Between):	-0.0002			
No. Observations:	333768	R-squared (Within):	5.197e-07			
Date:	Sun, Jul 29 2019	R-squared (Overall):	-8.903e-05			
Time:	14:55:35	Log-likelihood	-2.71e+05			
Cov. Estimator:	Unadjusted	F-statistic:	0.1733			
Entities:	264	P-value	0.6772			
Avg Obs:	1264.3	Distribution:	F(1, 333503)			
Min Obs:	8.0000					
Max Obs:	1277.0	F-statistic (robust):	0.1733			
		P-value	0.6772			
Time periods:	1277	Distribution:	F(1, 333503)			
Avg Obs:	261.37					
Min Obs:	255.00					
Max Obs:	263.00					
Parameter Estimates						
=====						
	Parameter	Std. Err.	T-stat	P-value	Lower CI	Upper CI
const	-1.037e-15	0.0009	-1.098e-12	1.0000	-0.0018	0.0018
Uber & Lyft	-0.0007	0.0017	-0.4163	0.6772	-0.0039	0.0026
=====						
F-test for Poolability: 2982.9						
P-value: 0.0000						
Distribution: F(263, 333503)						

OLS Regression Results						
=====						
Dep. Variable:	tickets	R-squared:	0.703			
Model:	OLS	Adj. R-squared:	0.703			
Method:	Least Squares	F-statistic:	2989.			
Date:	Sun, 21 Jul 2019	Prob (F-statistic):	0.00			
Time:	15:36:44	Log-Likelihood:	-2.7104e+05			
No. Observations:	333768	AIC:	5.426e+05			
Df Residuals:	333503	BIC:	5.454e+05			
Df Model:	264					
Covariance Type:	nonrobust					
=====						
	coef	std err	t	P> t	[0.025	0.975]
Intercept	-0.7892	0.015	-51.613	0.000	-0.819	-0.759
C(zone) [T. 1. 0]	0.2108	0.022	9.765	0.000	0.168	0.253
...
C(zone) [T. 263. 0]	0.0008	0.022	0.038	0.970	-0.042	0.043
Uber & Lyft	-0.0007	0.002	-0.416	0.677	-0.004	0.003
=====						
Omnibus:	79835.840	Durbin-Watson:	1.183			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	5088941.734			
Skew:	-0.080	Prob(JB):	0.00			
Kurtosis:	22.129	Cond. No.	265.			
=====						

FIGURE 4 Fixed Effects Model Result

1 DID: Difference in Differences

2 A pair of homogeneous taxi zones in **Figure 5** and **Figure 6** show how the parallel trend
 3 assumption was met. The map in **Figure 5** indicates that Flatiron and Midtown Center are similar.
 4 Most of their attributes, displayed in the table, are highly close. Their parking violation distribu-
 5 tion in 2014 is also similar, shown in blue in the first two line charts in **Figure 6**. However, the
 6 distribution of the number of Uber & Lyft trips and the parking tickets after 2014 of these two areas
 7 are different, shown by the orange and green lines on the line charts in Figure **Figure 6**. The Uber
 8 & Lyft trips in Midtown Center are almost twice the trips in Flatiron. Therefore, the difference
 9 between Uber & Lyft trips between the two homogeneous taxi zones is the only external factor that
 causes the significant change of the parking ticket.

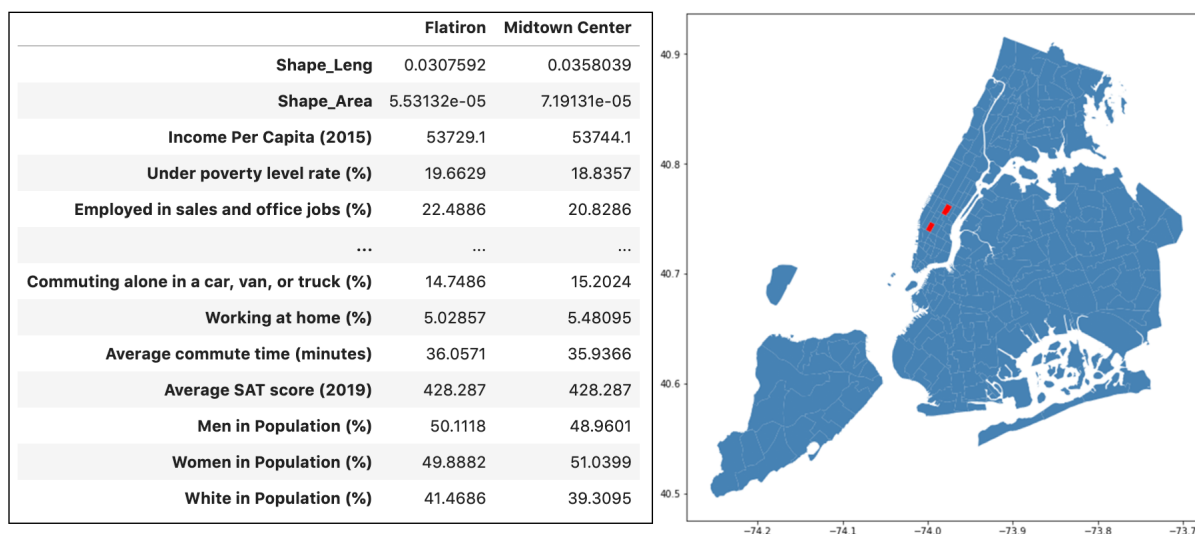


FIGURE 5 Features and Locations of a pair of highly homogeneous taxi zones

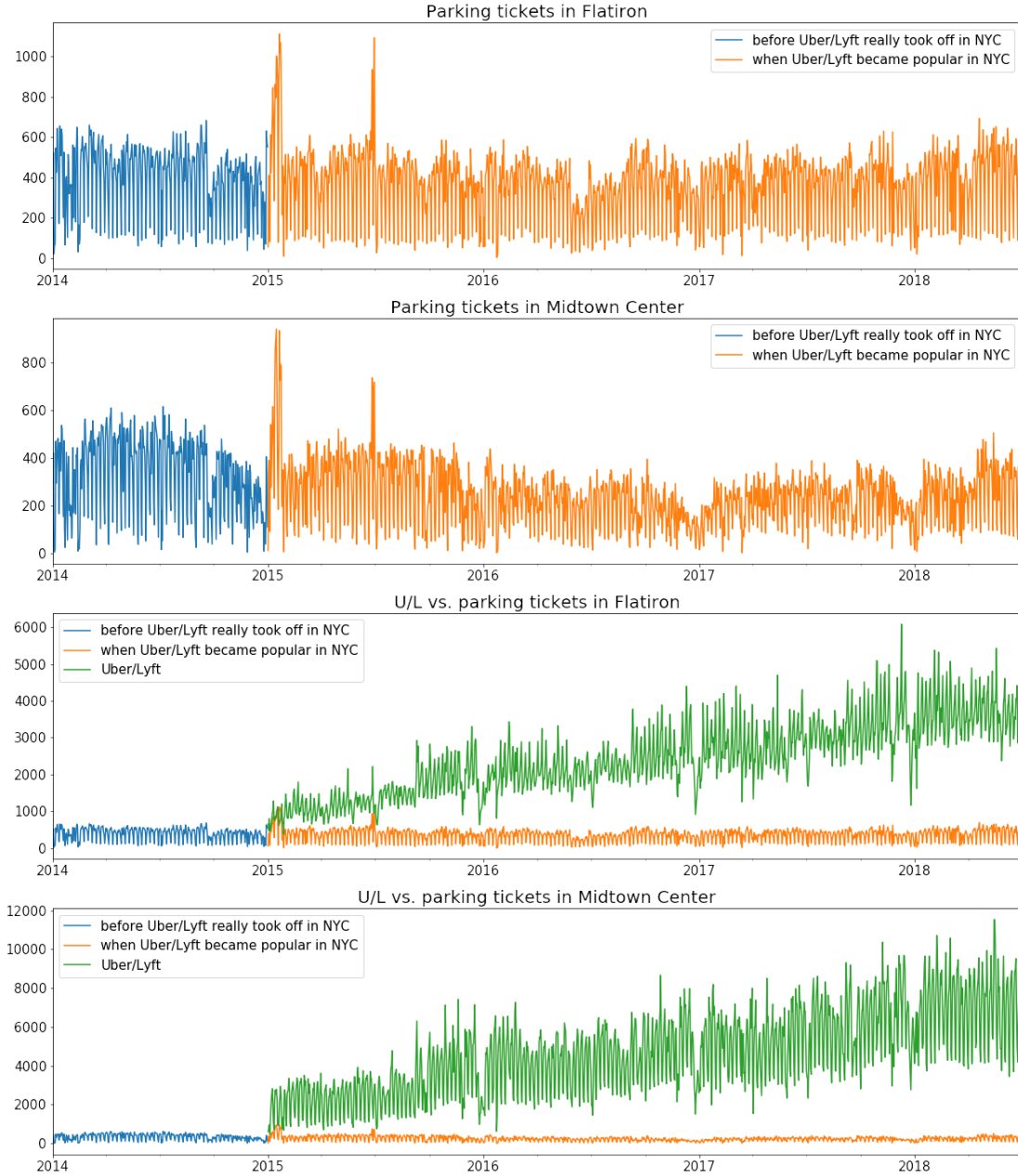


FIGURE 6 Tickets and TNC Trips of a pair of highly homogeneous taxi zones

1 Based on all qualified pairs, the DID model gave an average causal effect intensity of 224.2,
 2 which indicates that 1 parking ticket will be decreased with every 224.2 Uber or Lyft pickups
 3 increasing. So, the DID model suggests that Uber & Lyft rides can reduce parking violations in
 4 NYC and the effect is much more significantly than the result given by the previously explained
 5 FEM.

6 Pearson correlation analysis on the difference between Uber & Lyft trips and the corre-
 7 sponding difference between parking tickets, and hypothesis testing for non-correlation were used
 8 to validate the correctness of DID's conclusion. The results for all 900 qualified pairs among the
 9 whole NYC and grouped pairs within each borough are listed in the following **Table 3**:

	NYC	Manhattan	Brooklyn	Queens	Bronx	Staten Island
Coeff	0.021	-0.219	0.063	0.093	0.166	-0.034
P-value	0.534	0.028	0.612	0.420	0.225	0.908

TABLE 3 Pearson correlation analysis and hypothesis testing

1 The series of testing shows that although the DID concludes that increase of Uber & Lyft
2 leads to decrease of parking tickets in NYC, the correlation analysis, on the other hand, indicates
3 number of Uber & Lyft trips is positively related with number of parking ticket in overall NYC
4 and three specific boroughs, Brooklyn, Queens, and Bronx. However, it should be noted that
5 the statistical significance of this positive correlation is very low according to the very large P-
6 values. But at the same time, Manhattan's result has significantly supported the conclusion of
7 DID: negative correlation with p-value less than 0.05. Conversely, if we only use homogeneous
8 taxi pairs within the Manhattan, DID provides a even higher average estimate of casual intensity:
9 472. Based on the above results, we believe that at least in Manhattan, the conclusion that Uber &
10 Lyft have a negative causal effect on parking tickets is reliable.

11 CONCLUSIONS

12 The main goal of this paper was to explore the causal effect between ride-hailing services
13 represented by Uber & Lyft and parking violations: determining whether Uber & Lyft increases or
14 reduces the number of parking violations in NYC, and how strong the effect is. After a series of
15 data collection, preparation, and analyses, the results demonstrated that Uber & Lyft have negative
16 causal effect on NYC parking violations, suggesting Uber & Lyft are able to reduce parking vio-
17 lation. However, the extent of the impact has not been fully confirmed yet, due to the difference
18 between the result of the FEM and DID. Future improvements could be made to further modeling
19 the scale of the effect.

20 With more than 100 million users in the world (16), Uber & Lyft has far-reaching impacts;
21 some yet to be identified and explored. This paper serves as an example to apply large scale
22 machine learning methods on the available data to study and prove unnoticed yet potential effects
23 of ride-share companies on cities. Reducing parking violation can be regarded as a positive impact,
24 and also have following effects such as the reduction of parking enforcement pressure for the city.
25 In this aspect, Uber & Lyft could be beneficial to NYC, and potentially to other major cities in the
26 world. However, this is only one aspect of Uber & Lyft's influence. More research could be done
27 in this area to further explore other aspects of Uber & Lyft influence in the city.

1 **ACKNOWLEDGEMENTS**

2 This research was supported by Center for New York University's Urban Science + Process
3 and Wagner Graduate School of Public Service.

4 **AUTHOR CONTRIBUTIONS**

5 The authors confirm contribution to the paper as follows: Junjie Cai prepared datasets, as-
6 sisted in organizing the codes, and estimated Fixed Effect Model and Bayesian Network Model;
7 Junru Lu developed Difference in Difference Model, assisted in organizing the codes, and drafted
8 this manuscript; Pranay Anchan prepared data (Geocode FHV data with google API) and proofread
9 the manuscript; Shijia Gu helped in data collection and preparation, and model development. Yux-
10 uan Wang helped in data collection, model exploration, built the codes repository and proofread
11 the manuscript. Zhan Guo developed the research questions, provided feedback on model devel-
12 opment and policy implications, and proofread the manuscript. All authors reviewed the results
13 and approved the final version of the manuscript.

1 REFERENCES

- 2 1. Richter, W., *Uber and Lyft are gaining even more market share over taxis and rentals*,
3 2018.
- 4 2. Grocer, S., *How Uber and Lyft Compare, in Four Charts*, 2019.
- 5 3. statista, *Ride Hailing*, 2019.
- 6 4. Reed, E., *How Much Do Uber and Lyft Drivers Make in 2018?.*, 2018.
- 7 5. Brinklow, A., *Lyft, Uber Increase Traffic 180 Percent in Major Cities, Says Report*, 2018.
- 8 6. Schmitt, A., *Study: Uber and Lyft Caused U.S. Transit Decline*. *Streetsblog*, 2019.
- 9 7. Henao, A. and W. Marshall, A Framework for Understanding the impacts of ridesourcing
10 on transportation. In *Disrupting Mobility*, Springer, 2017, pp. 197–209.
- 11 8. Ackerman, S. S. and R. E. Moustafa, *Red Zone, Blue Zone: Discovering Parking Ticket
12 Trends in New York City*, 2011.
- 13 9. Pearl, J., Causal inference in statistics: An overview. *Statist. Surv.* 3, 2009, pp. 96–146.
- 14 10. Heckerman, D. and J. S. Breese, Causal Independence for Probability Assessment and In-
15 ference Using Bayesian Networks. *IEEE Transactions on Systems, Man, and Cybernetics*
16 *âÃ Part A: Systems and Humans*, Vol. 26, 1996, pp. 826–831.
- 17 11. Kim, J. H. and J. Pearl, A computational model for causal and diagnostic reasoning in
18 inference engines. *Proceedings /JCAf-83*, 1983, pp. 190–193.
- 19 12. Glenn Firebaugh, C. W. and M. Massoglia, *Handbook of Causal Analysis for Social Re-
20 search*, Springer, chap. 7, p. 113, 2013.
- 21 13. Angrist, J. D. and J.-S. Pischke, *Mostly harmless econometrics: An empiricist's compan-
22 ion*, Princeton University Press, Vol. 4, chap. 5, pp. 227–243, 2008.
- 23 14. Abadie, A., A. Diamond, and J. Hainmueller, Synthetic control methods for comparative
24 case studies: Estimating the effect of California's tobacco control program. *Journal of
25 the American statistical Association*, Vol. 105, No. 490, 2010, pp. 493–505.
- 26 15. Imai, K. and I. S. Kim, When Should We Use Unit Fixed Effects Regression Models for
27 Causal Inference with Longitudinal Data. *American Journal of Political Science*, Vol. 63,
28 No. 2, 2019, pp. 467–490.
- 29 16. Niu, E., *Uber Has Nearly 5 Times More Users Than Lyft*, 2019.