1    **SPATIO-TEMPORAL PATTERN ANALYSIS OF TAXI TRIPS IN NEW YORK CITY**
2
3
4
5    **Hartwig H. Hochmair**
6    University of Florida, Geomatics Program
7    3205 College Avenue
8    Ft. Lauderdale, FL-33314
9    Tel: 954 577-6317 Fax: 954 475-4125; E-mail: hhhochmair@ufl.edu
10
11
12
13
14
15
16
17   Word Count: 4977 words text +  8 tables/figures x 250 words (each)  = 6977 words
18
19   Date: November 15, 2015
20

1      **ABSTRACT**
2      A growing number of extensive datasets provides transportation planners with the necessary
3      means to analyze urban travel patterns and gain insight into urban dynamics. This paper explores
4      the spatial and temporal variation of taxi trips in New York City (NYC) by analyzing 29 million
5      trip records from a freely available dataset. The study examines the role of airports in trip
6      generation and attraction and its temporal changes, as well as the variation of travel speed during
7      the day for different districts within the study area. Comparison of hourly trip frequencies
8      between weekday and weekend in each district reveals similarities and differences in functional
9      drivers of taxi trip demand. A negative binomial regression model is presented which predicts
10     the number of taxi trips per district from subway, train, and bus infrastructure, socio-economic
11     and land use variables. As part of this approach eigenvector spatial filtering is applied to
12     explicitly model spatial autocorrelation in order to receive unbiased estimates and correct
13     inference. Independent of predictor variables, a combination of subway ridership and taxi trip
14     numbers for each district in a modemix variable allows, using Local Indicators of Spatial
15     Association statistics in a Geographic Information System (GIS), to identify districts that exhibit
16     an increased inclination towards taxi use and are currently poorly served by public transit. The
17     presented approach could be used as a decision support tool for deciding where investments in
18     rapid transit infrastructure and service would be particularly effective for increasing transit mode
19     share.
20
21
22     *Keywords*: travel demand, taxi trips, subway ridership, open data, GIS
23

1    **INTRODUCTION**
2    Taxi service can be viewed as a higher level of public transit by providing personal, on demand,
3    point-to-point transportation. It is used by travelers whose value of time and requirements on
4    comfort are higher than public users, or who are in need of such a service, such as (a) elderly,
5    young or low income people without their own vehicle who live in areas that are poorly served
6    by public transit, (b) people with disabilities needing door-to-door service, (c) tourists and
7    visitors who choose not to drive in an unfamiliar city or to use its public transit system (*1*), or (d)
8    people not owning a car undertaking a trip which cannot be easily done by transit, e.g. when
9    buying heavy grocery. This study explores spatial and temporal variations in taxi trip demand in
10   New York City based on a freely available set of taxi trip data, and extends methods and results
11   from of related studies that were conducted using similar taxi ridership datasets. A novel aspect
12   of this paper is the application of a spatial negative binomial regression model using eigenvector
13   spatial filtering (ESF) to explore the relationship between taxi trips and explanatory variables,
14   such a transit infrastructure, while explicitly considering spatial autocorrelation in observations.
15   Another novel aspect is the application of local statistic techniques in a combined variable of taxi
16   trip counts and subway boarding numbers, which can help to identify areas that are currently
17   poorly served by public transit. Furthermore, the paper will explore the suitability of the used
18   free taxi dataset for analyzing change of travel speed in subregions of the study area and its
19   change throughout the day..
20         Data used in the study stems from the Taxi and Limousine Commission of New York
21   City (TLC) and contains information about all medallion taxi trips in NYC in 2013. This freely
22   available data is an example of a growing number extensive datasets that provide transportation
23   planners with the opportunity to better understand urban travel patterns and urban dynamics.
24   Open data is based on the idea that certain data should be freely available to everyone to use and
25   republish, without restrictions from copyright or patents (*2*). NYC provides currently access to
26   over 1350 data sets on the Web (https://data.cityofnewyork.us/dashboard) in a variety of
27   machine-readable formats. The New York City Council requires all city agencies to open their
28   data by 2018. The taxi dataset used in this study has not yet been provided as open data by the
29   city, but it has been listed under planned releases on the city Website.
30
31   **PREVIOUS WORK**
32   Taxicab is a significant transportation mode in urban areas since it complements other public
33   transport modes with flexible door-to-door service through uninterrupted service (*3, 4*).
34   Therefore a thorough understanding of taxi trip behavior and its relationship to transit
35   infrastructure is important to further increase the efficiency of urban transportation systems.
36   Several recent studies focus on improving taxi efficiency. A study conducted in Taipei City
37   showed that 60–73% of their operation hours, taxi drivers were driving without passengers
38   because they did not know where potential customers were, leaving them no other choice than
39   wandering around the city (*5*). This problem was addressed in the paper by mining historical data
40   to predict passenger demand distributions with respect to contexts of time, weather, and taxi
41   location, using clustering algorithms. Another study applied time series forecasting techniques
42   on data from taxis equipped with real-time vehicle location systems to make short term
43   predictions on the passenger demand for 63 taxi stands in the city of Porto, Portugal (*6*). A
44   predictive model for the number of vacant taxis in a given area based on time of the day, day of
45   the week, and weather condition, using 2.6 million anonymous locations of 150 taxis in Lisbon,
46   Portugal, is presented in (*7*).

1    A wide range of spatial data sources have recently emerged due to advances in sensor and
2    telecommunication technology, informatics, and data processing, which facilitates the analysis of
3    various aspects of travel behavior. These include among others GPS tracking systems, which are
4    used for various transportation modes including cycling, walking, bus, and taxi. GPS tracking is
5    also used in all New York City taxis, which are regulated by the TLC. Over the past few years,
6    the TLC provided access to its trip database containing millions of taxi trips with temporal and
7    spatial information acquired by GPS, which includes taxi pickup and drop-off date, time, and
8    location, as well as fare and distance traveled.
9    Various recent research studies used this data source to analyze different aspects of taxi
10   ridership in NYC. One study analyzed travel times and found that travel times from truck-GPS
11   data can be better estimated from taxi-GPS data during AM and PM periods than during night
12   time, which indicates that speed differences between taxis and trucks are greater for free-flow
13   conditions (*8*). Another study used 10 months NYC taxi trip data from 2010 to estimate a
14   multiple linear regression model for each hour of the day to model pickups and drop-offs (*9*).
15   The results identified six important explanatory variables for taxi trips, which include population,
16   education, age, income, transit access time, and employment, where the influence of these factors
17   on taxi pickups and drop-offs changed at different times of the day. A related study extracted
18   one-week NYC taxi trip data from 2009 and explained trip frequency through six variables in
19   Ordinary Least Square (OLS) and geographically weighted regression (GWR) models (*4*).
20   Significant variables included the average commuting time, the density of people having a
21   bachelor's degree or higher, the median income level, the proportion of commercial area, the
22   road density and subway accessibility. The GWR model demonstrated a strong spatial variability
23   for parameter estimations and outperformed the OLS model in terms of model fit and mitigation
24   of spatial autocorrelation in residuals. The effect of winter storms on taxi trip characteristics was
25   explored in another study, showing that travelers opted to take shorter trips during the snowstorm
26   than during fair weather conditions (*10*). Taxi services were found to respond slowly to the
27   recovery in demand after the winter storm, possibly due to uncertainty over road conditions after
28   such an event.
29   Another study discusses the anticipated impacts of two recently enacted taxi regulation
30   changes on revenue increase, which are a fare increase, and the use of smart phone taxi
31   applications (*11*). The paper concludes that the fare increase is not expected to increase driver
32   incomes equally for all periods, and also that the apps' effectiveness is highly dependent on the
33   time of day and weather conditions. Using a classification and regression tree model on average
34   trip speed per mile on observed taxi trips for NYC travel time reliability was analyzed (*12*),
35   showing that reliability-related Day-of-Week and Time-of-Day periods did not follow the
36   conventional understanding of morning and evening, peak, and midday travel time patterns.
37   Using the notion of a shareability network, observed taxi trips were used to model the collective
38   benefits of sharing as a function of passenger inconvenience, and to efficiently compute optimal
39   sharing strategies on massive datasets (*13*). This simulation study has, however, not yet been
40   tested with real passengers, hence it not yet clear what constitutes minimal acceptable passenger
41   discomfort imposed by sharing a cab and caused delays.
42   Other analysis tasks that utilized NYC taxi trip data provided by the TLS include a binary
43   logit model to model the mode choice between transit and taxi mode (*14*), comparison of trip
44   characteristics between summer (July) and non-summer (March) months (*15*), and the
45   development of TaxiVis, which an analysis environment that implements a new model that

1 allows users to visually query taxi trips, under consideration of spatial, temporal, and attribute
2 constraints (*16*).
3
4 **STUDY SETUP**
5
6 **Study area and data sources**
7 The study area is New York city which covers five counties (Bronx, Kings, New York, Queens,
8 Richmond). Taxi trip data for September and October 2013 were downloaded from
9 http://www.andresmh.com/nyctaxitrips. These two months were chosen since they have no major
10 holidays and also had no unusual weather events reported. The data on the download Website
11 were originally obtained from the NYC Taxi and Limousine Commission through an FOIA
12 (Freedom of Information Act) request. Each trip record contains pickup and drop-off time and
13 location, passenger count, trip time and distance, and the driver's taxi permit license.
14        Average weekday and weekend ridership data for subway stations for the years 2009 to
15 2014 is available from the Metropolitan Transportation Authority (MTA) Website
16 (http://www.mta.info/). This ridership data together with the subway station geometry is
17 available in ESRI's Geodatabase format for download (*17*). The geodatabase contains also the
18 station geometries for three commuter rails (Long Island Rail Road, Metro-North Railroad, and
19 Staten Island Railway) within NYC boundaries as well as a line feature layer with the symbology
20 and geometry of the NYC subway system.
21        Population and household data at the census tract level were obtained from American
22 Community Survey 5-Year Estimates from 2009-2013 at the census tract level through ESRI's
23 Business Analyst. Longitudinal Employer-Household Dynamics (LEHD) Workplace Area
24 Characteristic data at the census block level was used to estimate the number of jobs per census
25 tract. Landuse (zoning) vector data was downloaded from MapPLUTO, which is hosted on the
26 NYC Department of City Planning Website.
27
28 **Data preparation**
29 The study area contains 2166 census tracts. To increase the number of taxi trips per analysis area
30 unit the census tract geometries were grouped into 130 districts using K-means clustering in
31 CrimeStat IV software. Socio-economic data (e.g., population, jobs, income, car ownership)
32 were aggregated accordingly. The taxi trip records as well as the district geometries were
33 imported into a PostgreSQL database with PostGIS extension The dataset for September and
34 October contained over 29.1 mio. trips. After importing trip records into the PostgreSQL
35 database following trips were removed:
36
37        •        missing coordinates for pickup or drop-off location
38        •        pickup or drop-off location outside the city boundaries
39        •        trip speed above 100 km/h or below 5 km/h
40        •        trip distance < 300 m
41
42 which resulted in 28.2 mio. trips that were retained for further analysis.
43
44 **Analysis methods**
45 Data analysis is split into two general parts. The first part analyzes patterns of taxi trips, whereas
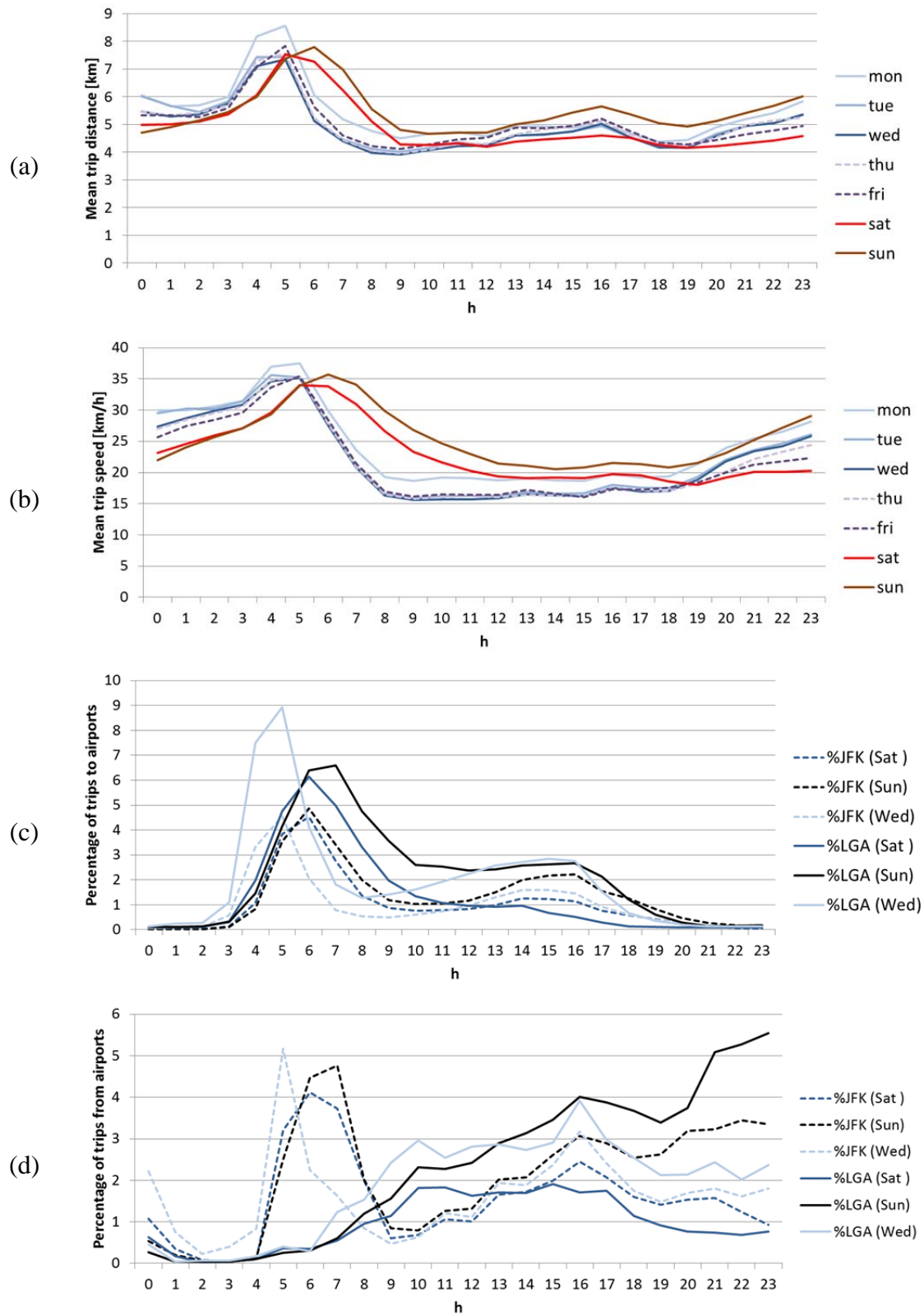46 the second part explores the relationship between taxi trips and other explanatory variables, such

1   as subway ridership. The first part is of exploratory nature and involves mapping and plotting
2   descriptive statistics related to trip frequency, trip distance, and trip speed, computed for
3   individual districts and/or time periods. It includes also an assessment of the similarity of hourly
4   trip demand patterns of weekdays and Sundays for districts, and uses local statistics (Anselin
5   Local Moran's I) to identify significant spatial clusters of districts with high or low similarity
6   values.
7          The second part uses multivariate regression at the district level to determine how transit
8   infrastructure (e.g., subway and bus stations) and subway ridership affect taxi trip demand after
9   controlling for non-network related variables, such as jobs and population. Since observed trip
10  numbers represent a count variable, a negative binomial regression model is used, which can
11  handle overdispersed count outcome variables. Spatial autocorrelation is handled through
12  eigenvector spatial filtering (ESF) A combined measure of taxi trip numbers and subway
13  ridership counts for each district will determine which district shows a tendency towards taxi
14  use. This variable can be analyzed through local statistics (Anselin Local Moran's I) as well to
15  identify local clusters of high taxi demand and thus areas where improved rapid transit
16  infrastructure could be particularly helpful to increase the transit mode share.
17
18  **PATTERN ANALYSIS OF TAXI TRIPS**
19
20  **Basic trip characteristics**
21  Typical recurring patterns of hourly and daily trip numbers of NYC cabs have been charted and
22  described elsewhere (*11*), revealing for example that the highest number of taxis on duty is
23  during weekend nights and lowest during Sunday evening. Furthermore, a drop in taxis available
24  can be observed during driver shift change, which takes place around 5am and 5pm. The
25  provided reference plots also the average distance per trip for 2010 taxi trip data, which is similar
26  to the curves plotted in FIGURE 1a and b for the analyzed two-month period in 2013. The
27  graphs indicate that mean trip distance and speed peak in the morning hours, with an earlier peak
28  during weekdays. The similar shape of the distance and speed curves indicate that longer trips
29  are faster, which can be expected given the higher proportion of limited access highways on
30  longer trips (*18*). There is no effect of rush hour traffic discernable for the overall system travel
31  speed, given that between 8am and 6pm weekday curves in both charts are relatively flat. This is
32  different when considering local trips only, i.e., after removal of long distance trips with their
33  higher average travel speed. A small portion of the morning peak in trip distances can be
34  attributed to the larger proportion of trips to airports (FIGURE 1c) or from airports (FIGURE 1d)
35  in the morning compared to other times of the day. Trips to and from the two major airports,
36  John F. Kennedy International Airport (JFK) and LaGuardia Airport (LGA), are longer on
37  average (JFK: ~25 km, LGA: ~15 km) than trips to and from other districts (~5 km). Another
38  contributing factor to longer trips in the morning could be the higher share of work related trips
39  (*19*).

1  **FIGURE 1 Trip distance (a), trip speed (b), and share of trip to airports (c) and from**
2  **airports (d) during a day.**

1
2   Mean trip distance varies strongly between analyzed districts. FIGURE 2 visualizes for
3   weekdays (Monday through Thursday) the mean travel distance (in km) for trips originating
4   from (FIGURE 2a) and ending in (FIGURE 2b) each mapped district, respectively. Only districts
5   with at least 100 pickups (FIGURE 2a) and drop-offs (FIGURE 2b) are considered. Taxi trip
6   distance decreases with distance from the CBD (Midtown and Lower Manhattan), indicating that
7   only a small portion of trips run between the CBD and more distant districts. In fact, the
8   percentage of weekday trips ending in three analyzed CBD districts and originating from JFK
9   among all trips ending in the CBD range only between 1.2% and 1.8%. These numbers are
10  somewhat higher for LGA with 2.5% to 2.9%, probably due to its closer proximity to the CBD
11  and its lack of access to the subway system. The maps show also that trips to and from airports
12  are longer than those to and from their surrounding districts.
13          The role of the CBD for taxi trips is more prominent when considering only trips that
14  begin or end at airports. FIGURE 2c visualizes for weekdays the spatial distribution of drop-offs
15  for trips originating from JFK, and FIGURE 2d shows where passengers were picked up for a
16  taxi ride to JFK. It is apparent that Midtown and lower Manhattan contribute strongly to taxi
17  demand from and to JFK. Furthermore the high pickup and drop-off numbers for the LGA
18  district shown in the maps indicate that taxi rides occur frequently between JFK and LGA. The
19  latter trips could stem from passengers who arrive from international flights at JFK and need to
20  connect to national flights departing from LGA, which does not serve international flights..
21  Results shown in FIGURE 2 are comparable to those for Sundays with only small differences.
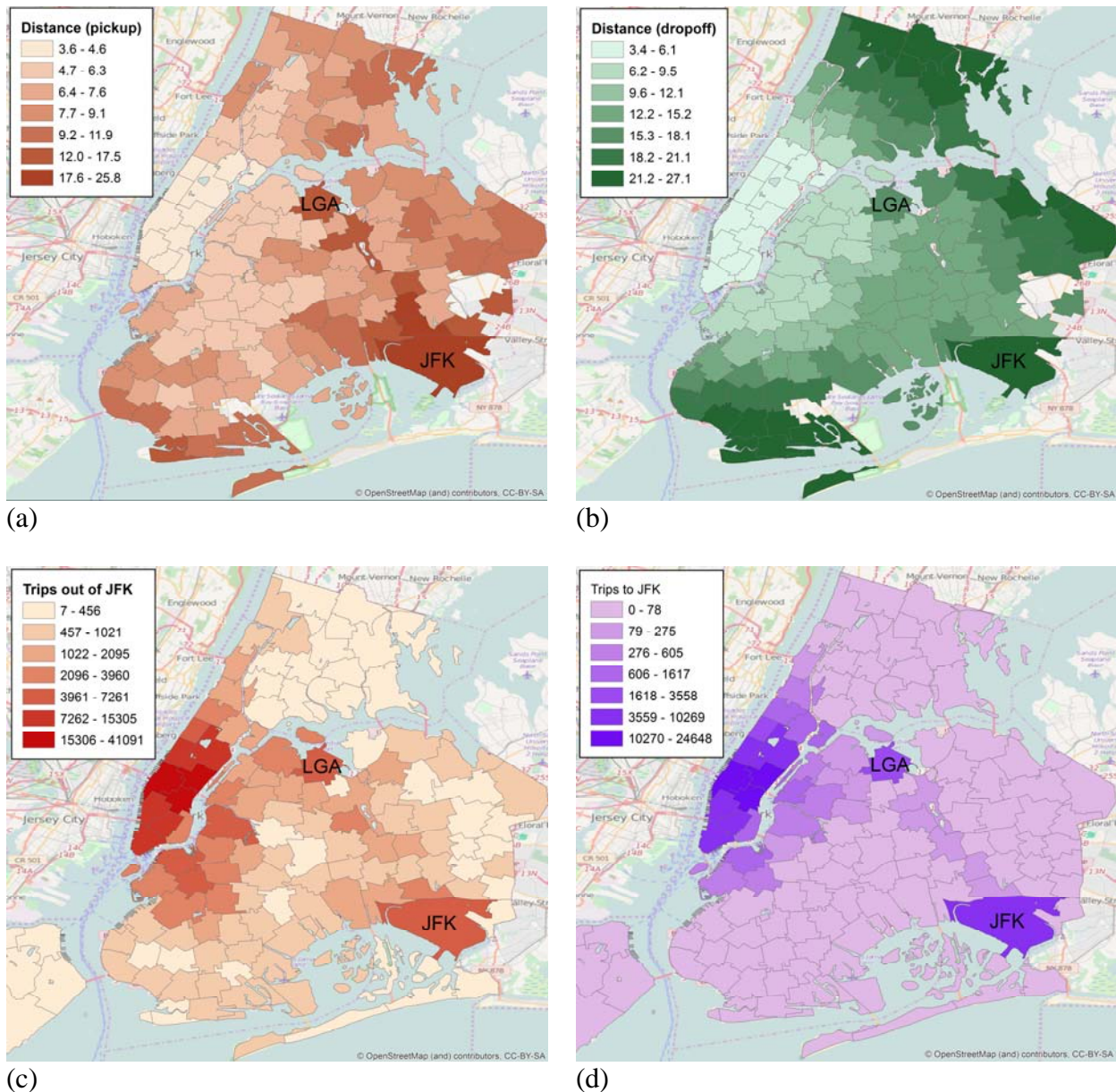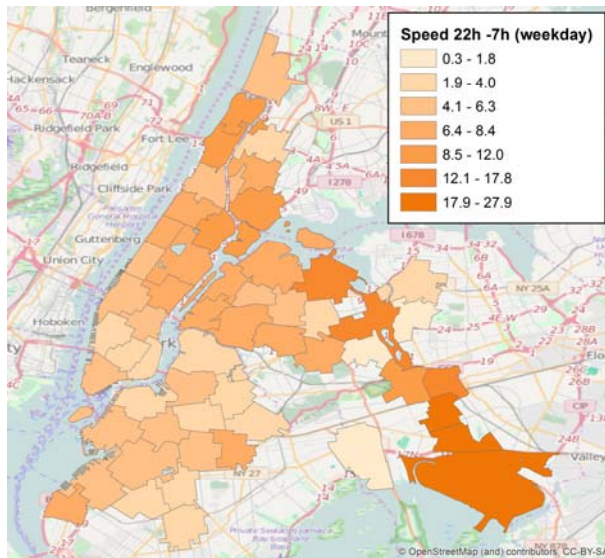22

(a)



(b)



(c)



(d)

**FIGURE 2 Weekday distance (in km) for trips originating (a) and ending (b) in mapped districts, and number of trip destinations (c) and origins (d) for trips from and to JFK.**
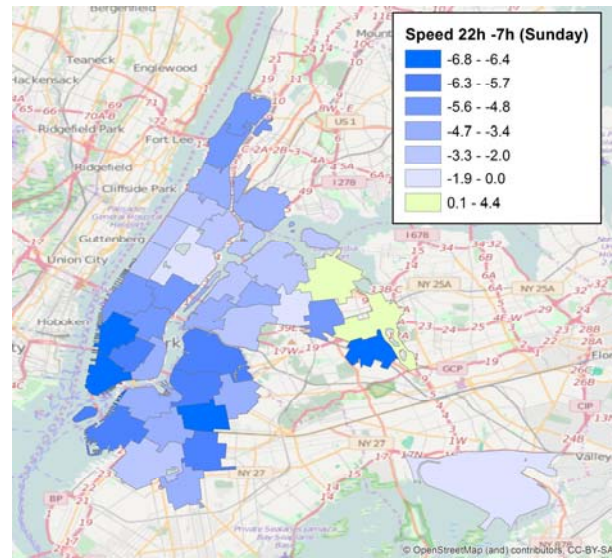
**Temporal variation of travel speed**

Historic or live floating car data are commonly used to determine urban road network travel time and congestion (*20, 21*). The goal of this analysis is to demonstrate that mean travel speed, as provided in the used dataset, can reproduce a pattern of changed travel speed over time as well. FIGURE 3 maps for different pairs of selected trip departure times the differences in mean travel speed (given in km/h). Only trips up to 8 km long are considered in order to focus on travel speed in and near each mapped district. Furthermore, only districts with at least 30 trip counts for both compared travel times are mapped (hence the gaps in the map). FIGURE 3a maps travel speed at 7am subtracted from that at 10pm for weekdays. The positive map values indicate that travel speed is in each district higher at 10pm than at 7am. The largest differences are found
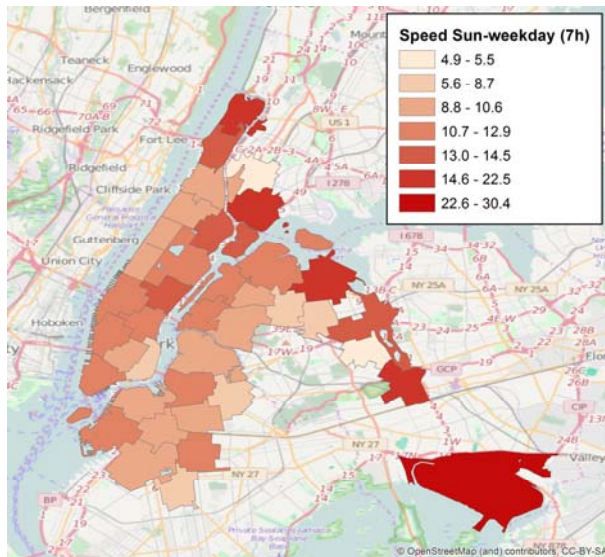
1    around the airports, where demand for taxi trips peaks in morning hours (compare FIGURE 1c
2    and d). In FIGURE 3b, which compares the travel speed between 7am and 10pm on Sundays,
3    most values are negative, suggesting faster travel in the morning than in the later evening.
4    FIGURE 3c compares travel speed at 7am between weekday and Sunday, indicating higher
5    travel speeds on Sunday morning. As opposed to this, FIGURE 3d shows that travel speeds are
6    comparable at 10pm between weekday and Sunday, since mapped differences are close to zero.
7

(a)                                                                          (b)

(c)                                                                          (d)

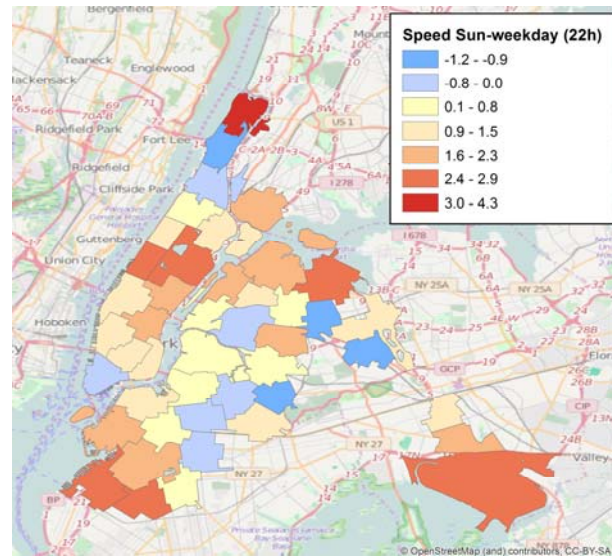8    **FIGURE 3 Difference in mean travel speed for different times of the day (a, b) and**
9    **different days (c, d). Numbers are given in km/h.**
10
11   **Weekend-weekday correlation in trip frequency**
12   The comparison of the hourly trip frequency pattern between weekend and weekday can give
13   insight into the functionality of neighborhoods (*22*). For each district with least 100 trips on

1  weekday and Sunday, the average number of trips was extracted for each hour on weekday and
2  Sunday. The obtained 24 trip counts for each day were then normalized between 0 and 1. The
3  similarity in taxi demand between weekday and Sunday was then determined through bivariate
4  correlation between the two normalized frequency series. The map in FIGURE 4a reveals
5  positive correlation values for both airports, indicating that the daily demand pattern for taxi trips
6  from the airport is similar throughout the week. Some districts in lower Manhattan show smaller
7  correlation values compared to their neighborhood, which could be because of nightlife activities
8  on weekends, or administrative buildings and offices (e.g. banks or city hall) which are closed on
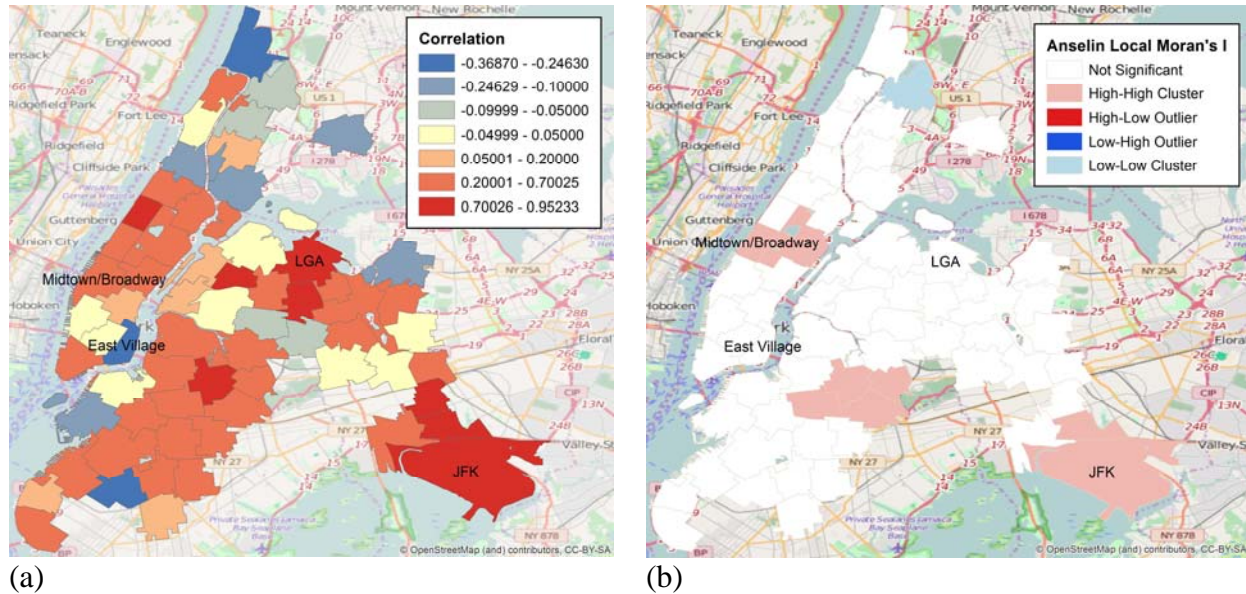9  weekends.
10



(a)                                                                              (b)

11  **FIGURE 4 Correlation in daily trip frequencies between weekday and Sunday (a), and**
12  **Anselin Local Moran's I (b).**
13

14  A local statistic is a descriptive statistic associated with a spatial dataset whose value varies from
15  place to place (*23*). Local Indicators of Spatial Association (LISA) statistics evaluate the
16  existence of clusters in a spatial arrangement of a variable (*24*). One example for LISA is
17  Anselin Local Moran's I, which can identify statistically significant spatial clusters of high or
18  low values and outliers. Applying local Moran's I to the autocorrelation variable in FIGURE 4a,
19  a High-High (HH) cluster identifies a district with a high correlation value surrounded by other
20  districts with high values, indicating a hot spot. HH clusters can be observed at JFK, around
21  Central park, and in a residential area in the north-east of Brooklyn (FIGURE 4b). A Low-Low
22  (LL) cluster indicates a district with a low correlation value that is surrounded by districts with
23  low values (cold spot). This occurs for one district in the north of the study area. The
24  computation of the levels of significance for clusters identified in FIGURE 4b (which are
25  visualized based on a threshold value of $p=0.05$) has not taken into account multiple testing and
26  spatial dependency. Depending on the correction method fewer clusters might potentially be
27  marked as significant.
28        The plots in FIGURE 5 show for four selected districts (labeled in FIGURE 4a) the
29  normalized trip counts per hour for weekday and Sunday, sorted by descending order of
30  correlation. JFK (FIGURE 5a) demonstrates the highest correlation among all districts, with a

1   morning peak around 6am for weekday and Sunday, possibly from arriving overnight flights.
2   LGA (FIGURE 5b) shows a steeper taxi demand between 6am and noon on weekdays than on
3   Sundays, possibly due to early arrivals of regional business flights during the week. The higher
4   trip demand between midnight and 6am for Sundays compared to weekdays in the
5   Midtown/Broadway district (FIGURE 5c) is probably due to the variety of late night
6   entertainment activities offered, which are primarily attended on weekend days (e.g., Saturday),
7   leading to an increase in trip frequency after midnight on the next day (i.e. Sunday). Taxi
8   demand picks up faster after 6am on weekdays compared to Sundays due to work trips in the
9   early morning. East Village (FIGURE 5d) is the district with the largest negative correlation. It
10  shows high traffic demand after midnight on Sunday (similar to the Midtown/Broadway district),
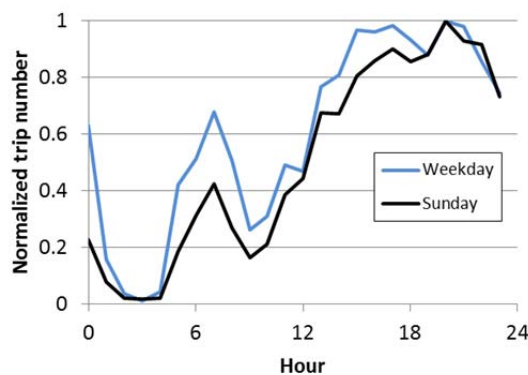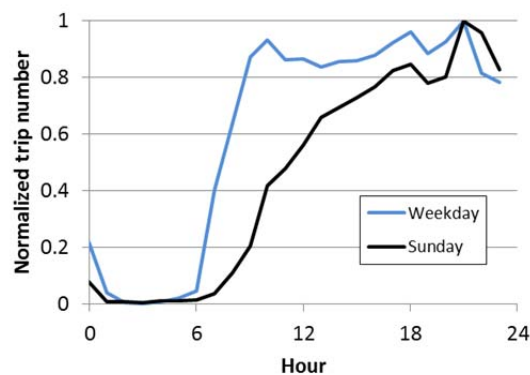11  but a lower relative demand afterwards compared to weekdays.
12

(a) JFK; r = 0.95                                   (b) LGA; r = 0.86

(c) Midtown/Broadway; r = 0.30                      (d) East Village; r = -0.37



13  **FIGURE 5 Correlation of daily trip frequencies between weekday and Sunday for selected**
14  **districts.**
15
16  **RELATIONSHIP BETWEEN TAXI TRIPS AND PUBLIC TRANSIT**
17
18  **Multivariate regression model**
19  TABLE 1 shows the results of selected (nonspatial) negative binomial regression models
20  (NBMs) for estimating the number of taxi trips from and to districts for weekdays (Monday
21  through Thursday), based on observations from all 130 districts. The models are arranged in a
22  descending order of AIC (Akaike information criterion), where a smaller AIC means a better
23  model fit. Due to bivariate correlations between predictor variables only a subset of predictor

1   variable can be included at a time, otherwise multicollinearity occurs. The sequence of models
2   tends to include coefficients from all areas considered. These include:
3
4       •       Total population and total number of jobs
5       •       Socio-economic variables: % of age group, % African American, % Hispanic,
6               % Bachelor of Science degree or higher, % of households without automobile
7       •       Land use (zoning): % manufacturing, % residential, % commerce, % parks, CBD
8       •       Build environment: Road density [km/km$^2$]
9       •       Transit Metro and train stations, bus stations, metro boardings
10
11  Models 1 through 5 (M1-M5) include the number of subway/train stations and bus stations as
12  predictor variables, expressing supply (TABLE 1). The subway/train stations variable is replaced
13  by subway ridership (a demand variable) by the variation of some models shown in the lower
14  part of TABLE 1. All five model estimations (and their variations in the lower half of the table)
15  have in common that population is positively associated with taxi trip counts and airports, which
16  could mean that more people (on airports or elsewhere) increase taxi demand, or that taxis target
17  areas with many potential customers (including airports) more frequently. In none of the models
18  CBD was significant. In M1 the number of subway/train stations is positively associated with
19  taxi trip counts. This is consistent with previous findings where a high correlation between
20  public transit ridership and taxi trips was explained by the direct demand for taxi service from
21  major transit stations (*1*). Furthermore, the model reveals an increased number of bus stops as
22  associated with a decrease in taxi ride numbers. This could indicate that bus trips in well served
23  areas compete with taxi transportation mode. Ethnicity variables became only significant in
24  combination land use variables (M2), but did generally not improve model fit. For all models
25  except for M1, the subway/train station variable becomes non-significant, indicating that much
26  of the variation of taxi ridership associated with rapid transit stations is absorbed by other land
27  user or socio-economic variables. M3 shows that an increase in the percentage of population
28  between 0 and 19 years is associated with a decrease in taxi use. A possible explanation is that
29  many families with children own an automobile and do not have to rely on taxi transport. M3
30  shows also that parks and residential areas are associated with fewer taxi trips, which could
31  indicate that taxi trips are often used for business trips in more commercial areas. M4 shows that
32  districts with a higher percentage of households without automobile tend to increase taxi trips, as
33  can be expected. Whereas for all other NBMs (M1, M2, M3, M5) the Variance inflation factor
34  (VIF) was < 3.2, indicating that multicollinearity among predictor variables did not pose a
35  problem in the model specifications (*25*), for M4 the VIF is 4.8, due to the high bivariate
36  correlation of the auto ownership variable with other variables, e.g. number of jobs. Model M5,
37  which includes also education and road density, was found to be the best fitting model while
38  keeping the VIF low. All models have in common that bus station is the only significant transit
39  factor associated with taxi trips. This finding can therefore also be considered by smaller
40  municipalities that do not have a subway or train system. A log-likelihood test in all identified
41  NBMs reported overdispersion in the count data at p <0.0001, indicating that a  NBM is better
42  suited in modeling taxi trip numbers than a Poisson regression model.
43          In the lower group of NBMs the job and income variables become non-significant,
44  meaning that subway boardings numbers absorb some of the effect of job numbers on taxi
45  demand. In the lower versions of M2 and M3 subway ridership remains a significant predictor,

1 whereas in M4 and M5 of the lower models it becomes non-significant (hence these modified
2 models are not shown in the table).
3
4 **TABLE 1 Negative Binomial Regression Estimates for weekday taxi trips from districts**

| Parameter | M1 | M2 | M3 | M4 | M5 |
|---|---|---|---|---|---|
| Intercept | **5.268\*\*** | **5.782\*\*** | **14.452\*\*** | **4.144\*\*** | **11.148\*\*** |
| Population | **6.313E-05\*\*** | **6.532E-05\*\*** | **7.541E-05\*\*** | **2.261E-05\*\*** | **6.272E-05\*\*** |
| % 0-19 years | | | **-0.213\*\*** | **-0.114\*\*** | **-0.254\*\*** |
| % Afr. American | | **-0.017\*\*** | | | |
| % BS | | | | | **0.051\*** |
| Jobs | **1.461E-05\*\*** | **1.040E-05\*\*** | **8.480E-06\*\*** | **7.528E-06\*** | **5.995E-06\*-** |
| Income | 1.684E-05$^+$ | | | **2.534E-05\*\*** | |
| Subway/train station | **0.137\*\*** | 0.058 | -0.035 | -0.058 | 0.017 |
| Bus stops | **-0.027\*\*** | **-0.026\*\*** | **-0.024\*\*** | **-0.013\*** | **-0.016\*\*** |
| Airport | **7.109\*\*** | **7.801\*\*** | **5.727\*\*** | **8.140\*\*** | **7.307\*\*** |
| % Residential | | **-0.044\*\*** | **-0.045\*\*** | | **-0.051\*\*** |
| % Park | | | **-0.060\*\*** | | **-0.042\*\*** |
| % HH without car | | | | **0.100\*\*** | |
| Road density | | **0.222\*\*** | | | **0.187\*\*** |
| AIC | 2265.6 | 2226.2 | 2196.4 | 2184.7 | 2172.4 |
| | | | | | |
| Intercept | **5.603\*\*** | **5.786\*\*** | **14.339\*\*** | | |
| Population | **5.184E-05\*\*** | **5.732E-05\*\*** | **6.646E-05\*\*** | | |
| % 0-19 years | | | **-0.208\*\*** | | |
| % Afr. American | | **-0.017\*\*** | | | |
| % BS | | | | | |
| Jobs | | | | | |
| Income | 1.491E-05 | | | | |
| Subway boardings | **2.731E-05\*\*** | **1.314E-05\*\*** | **8.638E-06\*** | | |
| Bus stops | **-0.024\*\*** | **-0.022\*\*** | **-0.022\*\*** | | |
| Airport | **7.179\*\*** | **7.735\*\*** | **5.806\*\*** | | |
| % Residential | | **-0.045\*\*** | **-0.045\*\*** | | |
| % Park | | | **-0.062\*\*** | | |
| Road density | | **0.228\*\*** | | | |
| AIC | 2262.5 | 2222.8 | 2197.0 | | |

5 Note: ** $p<0.01$, * $p<0.05$, $^+$ $p<0.1$ (statistical trend)
6
7 **Spatially filtered negative binomial model**
8 A fundamental assumption of regression analysis is residual independence. Incorporating spatial
9 data in nonspatial models typically results in residuals that exhibit spatial autocorrelation. Using
10 the regression results of the nonspatial prediction model for the weekday taxi trips reaching a
11 district as an example, residual diagnostics reveals significant autocorrelation with a Moran's I of
12 0.1120 (p=0.0087). To account for this spatial effect, a spatially filtered NBM is estimated which
13 uses eigenvector spatial filtering (ESF) to model spatial autocorrelation. A detailed description of

1 involved steps is provided in (*26, 27*). The spatial filter is used as a predictor variable in the
2 negative binomial regression model and expected to explain a considerable part of the variance
3 in the taxi trip crime distribution.
4
5 **TABLE 2 Estimation results for a nonspatial and spatially filtered negative binomial model**
6 **for weekday taxi trips to districts.**

| Parameter | Nonspatial Neg. Binomial Model | | | Spatially filtered Neg. Binomial Model | | |
|---|---|---|---|---|---|---|
| | Coefficient | Standard error | p-value | Coefficient | Standard error | p-value |
| Intercept | **11.273**** | **0.859** | **0.000** | **9.066**** | **0.586** | **0.000** |
| Population | **4.429E-05**** | **3.631E-06** | **0.000** | **1.815E-05**** | **2.771E-06** | **0.000** |
| % 0-19 years | **-0.184**** | **1.970E-02** | **0.000** | **-9.781E-02**** | **1.422E-02** | **0.000** |
| % BS | **3.259E-02*** | **1.306E-02** | **0.013** | 9.550E-04 | 8.932E-03 | 0.915 |
| Sub/train sta. | 3.317E-02 | 3.002E-02 | 0.269 | 2.278E-02 | 2.007E-02 | 0.256 |
| Bus stops | **-6.048E-03**** | **2.173E-03** | **0.005** | 2.655E-04 | 1.514E-03 | 0.861 |
| Airport | **4.465**** | **0.736** | **0.000** | **3.455**** | **0.499** | **0.000** |
| % Residential | **-2.922E-02**** | **5.435E-03** | **0.000** | **-8.287E-03*** | **3.894E-03** | **0.033** |
| % Park | **-2.014E-02*** | **8.098E-03** | **0.013** | 6.407E-04 | 5.481E-03 | 0.907 |
| Road density | **8.664E-02**** | **2.550E-02** | **0.000** | **6.666E-02**** | **1.717E-02** | **0.000** |
| Spatial filter | | | | **0.629**** | **4.2382E-02** | **0.000** |
| AIC | 2626.1 | | | 2504.2 | | |

7 Note: ** $p<0.01$, * $p<0.05$, $^+$ $p<0.1$ (statistical trend)
8
9 The results in TABLE 2 show that the spatial filter is highly significant ($p<0.0001$). It eliminates
10 all of the unexplained residual spatial autocorrelation, as found through a non-significant
11 Moran's I of -0.0085 ($p=0.9845$). This model satisfies therefore the fundamental assumption of
12 residual independence in count models (*28*). Comparison with the nonspatial NBM shows a
13 better model fit, as indicated by a lower AIC score. Model comparison shows that the spatially
14 filtered NBM reduces both the statistical significance as well as the magnitude of all coefficients
15 that were identified as significant in the nonspatial NBM, indicating that violating the
16 assumption of spatial independence in residuals leads to biased coefficients with inflated levels
17 of significance. Compared to the nonspatial NBM, in the spatially filtered NBM the number of
18 bus stops is no longer a significant predictor. Some land use variables (residential area, road
19 density) as well population (total and percent between 0 and 19 years) remain the only variables
20 with significant coefficients, besides the spatial filter. This means that based on global regression
21 models there is no clear direction of the effect of transit supply, e.g. subway, train and bus
22 stations, on taxi trip demand. Instead the direction of the effects, if present, may vary over space
23 and time, which can be explored through independent local spatial models (*4*), separate models
24 by the hour (*9*), or by adding interaction terms with location and time of day in the regression
25 model. In this regard, spatially-filtered NBMs are a reminder that results of nonspatial regression
26 models should be interpreted with care, especially when spatial autocorrelation is present in
27 residuals.
28
29 **Joint spatial pattern of taxi rides and subway use**
30 The following analysis aims to identify districts of the city with an increased tendency towards
31 taxi use compared to metro use, based on observed ridership data. Previous regression analysis

1   did not identify a clear relationship between subway supply and taxi demand (see TABLE 1,
2   TABLE 2). That is, presence of metro stations will in some cases increase the demand for taxi
3   trips, e.g. for last mile trips. A city may also have districts where taxi demand is high because of
4   a lack of rapid transit infrastructure and no other realistic transportation alternatives than taxi
5   (e.g. due to low car ownership of travelers). Such areas have a potential to increase rapid transit
6   ridership, e.g. through infrastructure improvements or policy changes. Although the NYC transit
7   network operates besides subway also three commuter rails, the following analysis considers
8   subway ridership only and is therefore limited to 101 districts that contain subway stations or are
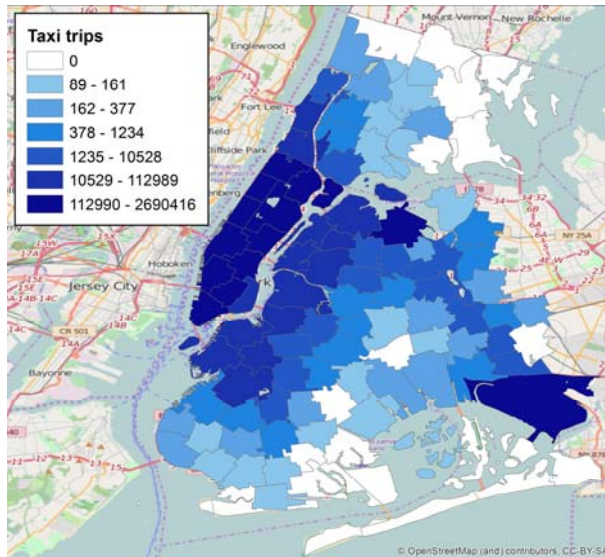9   located in-between them (FIGURE 6).
10         To describe the tendency of a district towards taxi or metro use numerically a modemix
11  variable *M* is introduced and computed for each district *d* as
12
13        $M_d = (Taxitrips_d \, / \, MaxTaxirides) - (Subwayridership_d \, / \, MaxSubwayridership)$            (1)
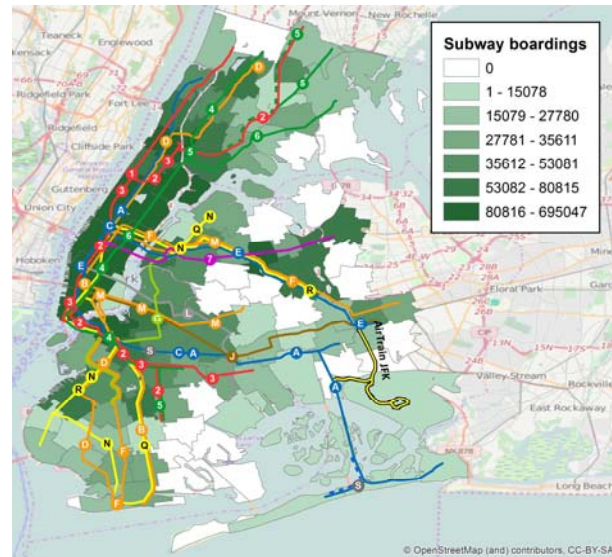14
15  where *Taxitrips* and *Subwayridership* are the total taxi trip count and subway ridership in district
16  *d* for a given time period, and *MaxTaxirides* and *MaxSubwayridership* are the maximum values
17  of *Taxitrips* and *Subwayridership*, respectively, in the whole study area. *M* ranges between -1
18  and 1. An *M* value below zero reveals a travel behavioral tendency towards subway use, and an
19  *M* value above zero indicates a tendency towards taxi use. Districts without access to subway but
20  observed taxi trips will always receive an *M* value above zero since the second term in the
21  equation is zero. Among this group of districts, those with a large number of taxi rides, e.g.
22  triggered by many jobs or a large population, will tend to have higher *M* values. For such
23  districts, providing rapid transit access could be an effective way to reduce the number of taxi
24  rides and substitute them with transit rides.
25         FIGURE 6a maps for each district the number of taxi trips that pick up passengers on a
26  weekday, and FIGURE 6b shows the total subway ridership on a weekday. JFK does not have a
27  subway station but an AirTrain system which connects to Jamaica and Howard Beach station.
28  Using the daily number of AirTrain trips, the known split of AirTrain users boarding at Jamaica
29  and Howard Beach station (approximately 60/40), and assuming a 50/50 split of JFK Train users
30  to continue their journey in the Long Island Road System and the subway system, respectively,
31  the corresponding number of daily AirTrain ridership was shifted from Jamaica and Howard
32  Beach station to the JFK analysis district. FIGURE 6c visualizes the modemix variable,
33  revealing higher values, i.e. a tendency towards taxi use, for the two airports and large portions
34  of Manhattan. FIGURE 6d maps the result of the Anselin Local Moran's I procedure on the
35  modemix variable (without considering multiple testing and spatial dependency). The High-Low
36  (HL) category indicates a high value district surrounded by low value districts, and the Low-
37  High (LH) category indicates the opposite. Both HL and LH cases indicate a spatial outlier. A
38  region of HH clusters, i.e. a hot spot of taxi demand, can be clearly discerned in the Manhattan
39  area. This high taxi demand relative to metro use could be caused by a limited capacity of the
40  subway system (e.g. crowded trains and stations) during rush hour. Another explanation could be
41  a reduced sensitivity of travelers in these districts to higher taxi trip fares (compared to subway),
42  given high household income (e.g. around Central Park), and well paid jobs in these districts. Yet
43  another possible reason is that the CBD has the lowest trip distance in NYC (see FIGURE 2a, b),
44  which keeps taxi fares generally affordable. Given that Manhattan has already a dense network
45  of subway lines, adding more subway infrastructure will most likely not reduce taxi ridership.
46  The situation is different for the marked High-Low outlier around the LGA airport where the
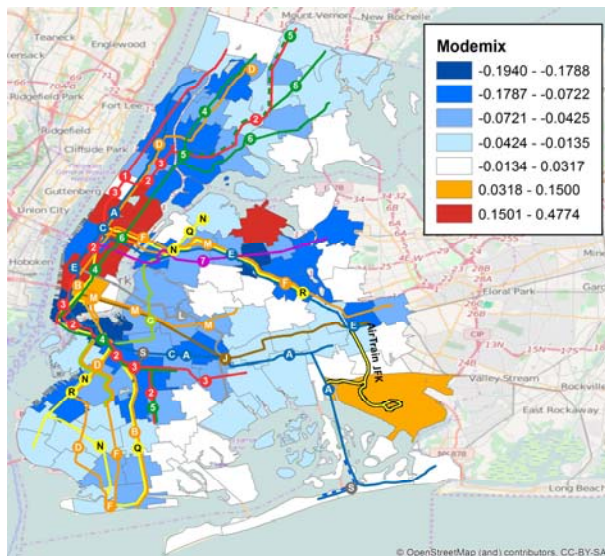
1   significantly higher modemix value of the LGA district compared to its surrounding districts
2   could be attributed to missing rapid transit infrastructure at LGA. As opposed to some of its
3   surrounding districts LaGuardia offers currently no subway or rail service. The only public
4   transportation provided is by several bus lines. Because of this deficiency, a new GLA AirTrain
5   system is currently proposed, which will be a 1.5-mile-long (2.4 km) people mover system
6   connecting with the subway and Long Island Rail Road in a similar manner to AirTrain JFK
7   (*29*). LISA results in FIGURE 6d support this endeavor by showing the outlier with an unusually
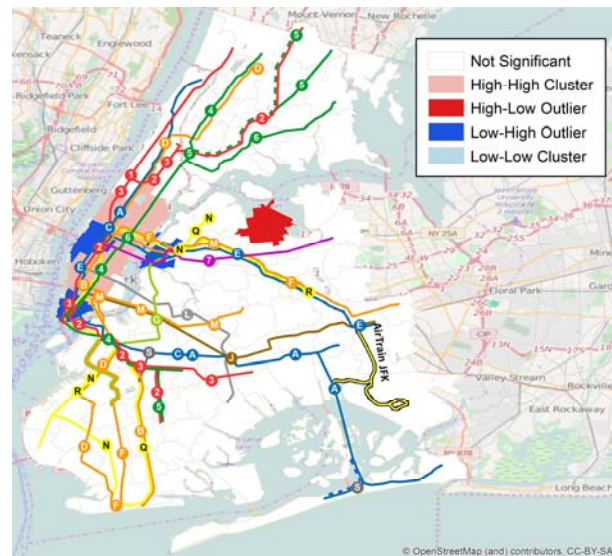8   high taxi mode share in the LGA area.
9



(a)

(b)

(c)

(d)

10  **FIGURE 6 Average number of taxi trips (a) and subway boardings (b) on a weekday, the**
11  **resulting modemix (c), and Anselin Local Moran's I (d).**
12

1    **SUMMARY AND FUTURE WORK**
2    This study explored various aspects of travel behavioral analysis in an urban environment using a
3    freely available dataset of taxi trip records for NYC, as well as subway ridership data at stations.
4    We expect additional extensive travel related datasets to become available for free in the near
5    future since numerous public agencies around the world are already actively participating in
6    open data initiatives. This trend will greatly enhance opportunities for transportation planners
7    and researchers to better understand behavioral patterns and urban dynamics. Aggregate analysis
8    for the complete study area showed a clear difference in taxi demand throughout the day when
9    comparing weekdays and Sundays patterns, which we attribute to different types of activities
10   conducted on weekdays and Sundays. A more refined picture of the variation of activities
11   between different days and regions can be obtained by correlating hourly trip count values on
12   weekdays and Sundays for each district. The analysis of ridership frequency throughout a day,
13   especially when used from multiple days, can be used to distinguish between different functional
14   areas, such as airports, residential areas, business districts, or entertainment districts. Combining
15   taxi trip counts with subway ridership data was shown to provide useful information about areas
16   that might be currently underserved by public transport, and where observed high taxi trip
17   numbers indicate a higher demand for alternative transportation, such as rapid transit. To draw
18   more information conclusions, the analysis would need to include heavy rail and bus services,
19   which is deferred to future work. Some of the nonspatial binomial regression models presented
20   reveal that subway and train stations are associated with an increase in taxi trips, confirming
21   results from previous studies (*1, 4*). This effect, becomes, however statistically nonsignificant
22   when using a spatially filtered binomial regression model that explicitly models spatial
23   autocorrelation in statistical analysis, showing that this relationship does not hold throughout the
24   study area, but may have to be addressed on local models. Bus stations tend to be associated with
25   a lower taxi demand. Also, this effect disappears in the global model when taking into account
26   spatial autocorrelation.
27
28   **REFERENCES**
29

30   (*1*)    Kattan, L., de Barros, A., and Wirasinghe, S. C. Analysis of work trips made by taxi in
31            Canadian cities. *Journal of Advanced Transportation*, Vol. *44*, 1, 2010, pp. 11-18.
32   (*2*)    Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., and Ives, Z. DBpedia: A
33            Nucleus for a Web of Open Data. In K. Aberer, K.-S. Choi, N. Noy et al. (Eds.), *The*
34            *Semantic Web*, *LNCS 4825*, Springer, Berlin, 2007, pp. 722-735.
35   (*3*)    King, D. A., Peters, J. R., and Daus, M. W. *Taxicabs for Improved Urban Mobility: Are*
36            *We Missing an Opportunity?* Transportation Research Board - 91st Annual Meeting,
37            Washington, D.C. Transportation Research Board of the National Academies, 2012.
38   (*4*)    Qian, X. and Ukkusuri, S. V. Spatial variation of the urban taxi ridership using GPS data.
39            *Applied Geography*, Vol. *59*, 2015, pp. 31-42.
40   (*5*)    Chang, H.-w., Tai, Y.-c., and Hsu, J. Y.-j. Context-aware taxi demand hotspots
41            prediction. *International Journal of Business Intelligence and Data Mining*, Vol. *5*, 1,
42            2010, pp. 3-18.
43   (*6*)    Moreira-Matias, L., Gama, J., Ferreira, M., and Damas, L. *A Predictive Model for the*
44            *Passenger Demand on a Taxi Network*. 15th International IEEE Conference on Intelligent
45            Transportation Systems Anchorage, Alaska, USA. IEEE, 2012.

(*7*)  Phithakkitnukoon, S., Veloso, M., Bento, C., Biderman, A., and Ratti, C. Taxi-Aware Map: Identifying and Predicting Vacant Taxis in the City. In B. d. Ruyter, R. Wichert, D. V. Keyson et al. (Eds.), *Ambient Intelligence*, *LNCS 6439*, Springer, Berlin, 2010, pp. 86-95.

(*8*)  Morgul, E. F., Ozbay, K., Iyer, S., and Holguín-Veras, J. *Commercial Vehicle Travel Time Estimation in Urban Networks using GPS Data from Multiple Sources*. Transportation Research Board - 92nd Annual Meeting, Washington, D.C. Transportation Research Board of the National Academies, 2013.

(*9*)  Yang, C. and Gonzales, E. J. Modeling Taxi Trip Demand by Time of Day in New York City. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. *2429*, 2014, pp. 110-120.

(*10*)  Qing, C., Parfenov, S., and Kim, L.-J. *Identifying Travel Patterns During Extreme Weather Using Taxi GPS Data*. Transportation Research Board - 94th Annual Meeting, Washington, D.C. Transportation Research Board of the National Academies, 2015.

(*11*)  Kamga, C., Yazici, M. A., and Singhal, A. Analysis of taxi demand and supply in New York City: implications of recent taxi regulations. *Transportation Planning and Technology*, Vol. *38*, 6, 2015, pp. 601-625.

(*12*)  Yazici, M. A., Kamga, C., and Mouskos, K. C. Analysis of Travel Time Reliability in New York City Based on Day-of-Week and Time-of-Day Periods. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. *2308*, 2011, pp. 83-95.

(*13*)  Santi, P., Resta, G., Szell, M., Sobolevsky, S., Strogatz, S. H., and Ratti, C. Quantifying the benefits of vehicle pooling with shareability networks. *PNAS*, Vol. *111*, 37, 2014, pp. 13290–13294.

(*14*)  Yang, C., Morgul, E. F., Gonzales, E. J., and Ozbay, K. Comparison of Mode Cost by Time of Day for Nondriving Airport Trips to and from New York City's Pennsylvania Station. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. *2449*, 2014, pp. 34-44.

(*15*)  Parfenov, S., Weeks, A., and Alam, Z. *Travel Patterns of NYC's Yellow Taxis: Routing, Activity and Results* ESRI International User Conference, San Diego, CA, 2014.

(*16*)  Ferreira, N., Poco, J., Vo, H. T., Freire, J., and Silva, C. T. Visual Exploration of Big Spatio-Temporal Urban Data: A Study of New York City Taxi Trips. *IEEE Transactions on Visualization and Computer Graphics*, Vol. *19*, 12, 2013, pp. 2149-2158.

(*17*)  Donnelly, F. Introduction to the NYC Geodatabase (nyc_gdb) ArcGIS Version. http://www.baruch.cuny.edu/geoportal/nyc_gdb/data/intro_nycgdb_arc.pdf. Accessed 7/27/2015

(*18*)  Parthasarathi, P., Hochmair, H. H., and Levinson, D. M. Network structure and spatial separation. *Environment and Planning B, Planning and Design*, Vol. *39*, 1, 2012, pp. 137-154.

(*19*)  Jiang, S., Ferreira, J., and González, M. C. Clustering daily patterns of human activities in the city. *Data Mining and Knowledge Discovery*, Vol. *25*, 3, 2012, pp. 478-510.

(*20*)  Li, Q., Zhang, T., and Yu, Y. Using cloud computing to process intensive floating car data for urban traffic surveillance. *International Journal of Geographical Information Science*, Vol. *25*, 8, 2011, pp. 1303-1322.

1    (*21*)    Rahmani, M., Jenelius, E., and Koutsopoulos, H. N. *Route Travel Time Estimation Using*
2             *Low-Frequency Floating Car Data* 16th Interntional IEEE Annual Conference on
3             Intelligent Transportation Systems, The Hague, The Netherlands, 2013.
4    (*22*)    Sagl, G., Delmelle, E., and Delmelle, E. Mapping collective human activity in an urban
5             environment based on mobile phone data. *Cartography and Geographic Information*
6             *Science*, Vol. *41*, 3, 2014, pp. 272-285.
7    (*23*)    O'Sullivan, D. and Unwin, D. J. *Geographic Information Analysis (2nd ed.)*, John Wiley
8             & Sons, Hoboken, New Jersey, 2010.
9    (*24*)    Anselin, L. Local Indicators of Spatial Association—LISA. *Geographical Analysis*, Vol.
10            *27*, 2, 1995, pp. 93-115.
11   (*25*)    Agresti, A. *Categorical Data Analysis (2nd ed.)*, John Wiley & Sons, Hoboken, NJ,
12            2002.
13   (*26*)    Chun, Y. and Griffith, D. A. *Spatial Statistics and Geostatistics*, SAGE Publications Ltd,
14            Los Angeles, 2013.
15   (*27*)    Helbich, M. and Arsanjani, J. J. Spatial eigenvector filtering for spatiotemporal crime
16            mapping and spatial crime analysis. *Cartography and Geographic Information Science*,
17            Vol. *42*, 2, 2014, pp. 134-148.
18   (*28*)    Griffith, D. and Haining, R. Beyond Mule Kicks: The Poisson Distribution in
19            Geographical Analysis. *Geographical Analysis*, Vol. *38*, 2006, pp. 123-139.
20   (*29*)    Alberts, H. R.  At Long Last, The AirTrain Will Go To LaGuardia Airport.
21            http://ny.curbed.com/archives/2015/01/20/at_long_last_the_airtrain_will_go_to_laguardi
22            a_airport.php. Accessed 7/30/2015