

1 **Using Location Based Social Network Data for Activity Intensity Analysis: A**

2 **Case Study of New York City**

3 September, 2019

4 **Haluk Laman**

5 Research Assistant

6 Department of Civil, Environmental & Construction Engineering

7 University of Central Florida

8 Tel: 1-407-414-4764

9 Email: [haluklaman@knights.ucf.edu](mailto:haluklaman@knights.ucf.edu)

10 ORCID number: 0000-0003-0884-610X

11

12 **Shamsunnahar Yasmin\***

13 Postdoctoral Associate

14 Department of Civil, Environmental & Construction Engineering

15 University of Central Florida

16 Tel: 407-823-4815, Fax: 407-823-3315

17 Email: [shamsunnahar.yasmin@ucf.edu](mailto:shamsunnahar.yasmin@ucf.edu)

18 ORCID number: 0000-0001-7856-5376

19

20 **Naveen Eluru**

21 Associate Professor

22 Department of Civil, Environmental and Construction Engineering

23 University of Central Florida

24 Tel: 1-407-823-4815, Fax: 1-407-823-3315

25 Email: [naveen.eluru@ucf.edu](mailto:naveen.eluru@ucf.edu)

26 ORCID number: 0000-0003-1221-4113

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19

**ABSTRACT**

Location-based social networks (LBSN) are social media sites where users’ check-in at venues and share content linked to their geo-locations. LBSN, considered as a novel data source, contain valuable information for urban planners and researchers. While earlier research efforts focused either on disaggregate patterns or aggregate analysis of social and temporal attributes, there is no attempt to relate the data to transportation planning outcomes. To that extent, the current study employs an LBSN service-based data for aggregate level transportation planning exercise by developing land-use planning models. Specifically, we employ check-in data aggregated at the census tract level to develop a quantitative model for activity intensity as a function of land use and built environments attributes for the New York City (NYC) region. A statistical exercise based on clustering of census tracts and negative binomial regression analyses are adopted to analyze the aggregated data. We demonstrate the implications of the estimated models by presenting the spatial aggregation profiling based on the model estimates. The findings provide insights on relative differences of activity engagements across the urban region. The proposed approach thus provides a complementary analysis tool to traditional transportation planning exercises.

*Keywords:* Location-based social networks, land use profiles, check-in, activity intensity, urban planning, clustering, negative binomial regression

# 1. INTRODUCTION

Smartphone ownership among Americans has rapidly risen to 77% in 2018 from 35% in 2011 (1). The ubiquity of smartphones with an embedded global position services (GPS) allows for obtaining precise individual level location information. In fact, according to a recent report by Pew Research Center (1), 90% of smartphone users obtain directions, recommendations and other location specific information from their phone. Several social networking sites (such as Twitter, Foursquare, Gowalla, and Facebook) allow users to share content on their websites with geo-coded information often referred to as location based social networks (LBSN). These location-based services allow users to “check-in” at a venue (such as restaurant or public park) based on their GPS coordinates providing them with location specific status update. While privacy concerns among users have ensured that the usage of location-based services is not universal, a large share of the population still adopts these services. For instance, 28% of American adults use a mobile or LBSN service. Furthermore, 12% of smartphone owners use their phone to check-in locations using the LBSN service. More interestingly, 7% of all adults allow the social media service they are using to automatically share their locations when they update their status. As is expected, the usage is higher among younger individuals - ages between 18-29 (16%), 30-49 (11%), 50-64 (9%) and 65+ (11%). These usage rates for LBSN clearly highlight the small share of adoption. However, given the large number of smartphone users, the data from these services would be larger than the data collected from traditional transportation data collection approaches (such as household surveys). Thus, it is not surprising that in recent years several studies have explored the use of such LBSN based datasets acquired from websites for data mining, land use planning, urban mobility analysis and transportation analysis (see (2)).

1           To be sure, the data available from LBSN services is not without limitations (as identified  
2 in (3)). First, the data does not provide detailed information (on gender, age, education) at the  
3 individual level. Second, even avid LBSN service users are unlikely to “check-in” every event in  
4 their day (particularly the routine activities). The activity start and end times are also unlikely to  
5 be available. Finally, the sample of individuals providing the information represents a sample of  
6 the population that is unlikely to be representative of the broader population. In fact, recently a  
7 study by Rzeszewski (4) illustrated that the data obtained from are inherently different based on  
8 the user behavior and the social media platform adopted. The authors cautioned analysts  
9 considering merging data from multiple platforms. Given these inherent biases, the data obtained  
10 from such services are prone to bias at a disaggregate level. On the other hand, employing the data  
11 for aggregate level analysis might provide a more representative population behavior. For  
12 example, rather than focusing on an individual’s activity locations, based on the LBSN data  
13 identifying number of “check-ins” at a spatial unit such as census tract might offer relative  
14 differences of activity engagement across the urban region. More importantly, the traditional data  
15 collection methods (such as household surveys) provide sparse information on such activity  
16 engagement information. Thus, employing LBSN data check-ins to identify activity centers (based  
17 on attractiveness) across the urban regions will provide a complementary analysis to traditional  
18 land use transportation planning exercises (5).

19           Given the large number of LBSN users, the data available provide us with large-scale  
20 datasets for activity analysis. The LBSN users provide analysts with detailed spatio-temporal data  
21 that can be utilized for planning applications. The main objective of our study effort is to employ  
22 LBSN service-based data for aggregate level planning exercise by developing land-use  
23 transportation planning models. To elaborate, using activity check-ins within a spatial aggregation

1 as a surrogate measure of attractiveness, the study provides a quantitative relationship between  
2 attractiveness and various socio-demographic, points of interest, transportation infrastructure and  
3 land-use attributes. The established relationship will allow transportation and land-use planners to  
4 identify what factors affect zonal/destination attractiveness and pro-actively plan for potential new  
5 demand with changing socio-demographic, land-use and/or transportation infrastructure patterns.  
6 Example of changes that can be analyzed include development of mixed use developments in dense  
7 neighborhoods or addition of public transit infrastructure.

8         In our research effort, data from LBSN provider Foursquare that allows users to check-in  
9 at indoor and outdoor venues (such as café, restaurant or public spaces) via smartphones is utilized.  
10 The geo-coded data is aggregated using Geographical Information System (GIS) techniques to  
11 obtain check-ins at a census tract level also referred to as “activity intensity” for the New York  
12 City (NYC) region. The relationship of the computed activity intensity variable with socio-  
13 demographics, land use variables, transportation and infrastructure variables, and points of  
14 interests at the census tract level is analyzed to offer insights on interconnectedness of activity  
15 intensity and other attributes.

16         A statistical exercise based on clustering and negative binomial regression analysis are  
17 adopted to analyze the aggregated data. The cluster analysis is performed in order to categorize  
18 the census tracts in the NYC region as a function of various exogenous variables. The clustering  
19 approach, rather than considering the entire city as homogenous allows us to distinguish across  
20 different clusters. Subsequently, cluster specific regression analysis is employed to identify the  
21 factors that affect the “check-ins” in the cluster. As the “check-ins” are non-negative integer  
22 values, negative binomial regression models were adopted for cluster specific regression models.  
23 The models estimated are employed to illustrate the impact of various parameters on check-ins

1 using a hot spot analysis. Hence, we illustrate how spatial distribution of activity patterns derived  
2 by LBSN data can be utilized to reveal urban activity patterns.

3 The remainder of the paper is organized as follows: Section 2, provides a review of earlier  
4 research and positions the current work in context. In Section 3, data source and description are  
5 provided. The research methods employed, model results, validation statistics and hot spot analysis  
6 are presented in Section 4. Finally, Section 5 concluded the paper.

7

## 8 **2. EARLIER RESEARCH AND CURRENT STUDY IN CONTEXT**

9 The traditional research efforts examine activity travel patterns (and related choices) based on  
10 traditional household travel surveys. The literature in this context is quite vast and it is beyond the  
11 scope of the paper to document (see (6) and (7) for a detailed summary of earlier work). With the  
12 increasing adoption of mobile devices, there is growing research employing innovative data  
13 sources for transportation planning analysis. In this context, we present the review of earlier studies  
14 along two streams: (1) research employing social media data for non-transportation research  
15 context and (2) research employing social media data for transportation research contexts.

16 The *first stream of studies* has mainly originated in the fields of social sciences and  
17 computer science. The emphasis of these research efforts is to extract behavioral insights on online  
18 activity and offline interactions (8). Cheng et al. (9) derived an algorithm to understand the  
19 mobility patterns of LBSN users. By studying several different metropolitan areas, user  
20 displacement, radius of gyration, and returning probability of individuals were determined. Their  
21 findings can be summarized as; LBSN users follow simple reproducible patterns which refer to  
22 Levy Flight type patterns, social status is coupled to mobility, and content analysis can reveal  
23 hidden context between people and locations. More recently, Ahas et al. (10) and Cao et al. (11)

1 employed mobile and/or location based social network data to study temporal and spatial  
2 differences in urban regions from multiple countries. A hierarchical statistical approach - the  
3 nested Chinese Restaurant Franchise (nCRF) - based on tweet contents of LBSN data of the US  
4 was proposed by Ahmad et al. (12) to infer a latent distribution of user locations. Kling and  
5 Pozdnoukhou (13) also used topic modeling for investigating space-time dynamics of time-  
6 stamped and geo-located check-in information. Topic modeling was also employed to process  
7 urban activity patterns and classify them for LBSN data of NYC (3).

8       The second stream of research has explored the viability of social media data for  
9 transportation planning purposes. Frias-Martinez et al. (5) used an unsupervised Neural Networks  
10 technique named Self Organizing Maps (SOM) to Manhattan area of NYC. The study findings  
11 indicate that LBSN data can serve as a complimentary source of information for urban planning  
12 development. Cranshaw et al. (9) employed the fine spatial resolution based on geo-located tweets  
13 by clustering nearby locations with similar activities and revealing social-spatial divisions in  
14 Pittsburgh. Wakamiya et al. (14) used LBSN data of three cities in Japan (Osaka, Nagoya, Tokyo)  
15 to study the crowd and individual movements across geography by aggregate and dispersion  
16 models as well as semantics of the tweet contents. They combined temporal analysis with k-means  
17 clustering based on the spatial check-ins and urban types by tracking common patterns in different  
18 regions. Their findings confirmed that crowd activities determined via Twitter can characterize  
19 living spaces in cities. Noulas et al. (15) used rank based movement model by ranking transitions  
20 by distance in order to capture urban mobility pattern variations. A similar study using rank-based  
21 models was conducted aiming to determine how a large-scale geo-location data set can be analyzed  
22 to classify and refer to individual activity patterns (11). As a result of aggregate and disaggregate  
23 level analysis in New York, Chicago and Los Angeles areas, the study concluded that people

1 choose their destinations mainly based on the popularity of these places. Bawa-Cavia (16)  
2 conducted inter urban analysis of Foursquare data in three metropolitan cities (NYC, London, and  
3 Paris) to understand difference in spatial structures across these cities. Zhan et al. (17) deployed  
4 supervised (random forest algorithm) and unsupervised (k-means clustering) approaches to infer  
5 land use of NYC based on LBSN data. The findings confirm that LBSN data can be used as a  
6 complementary data source in land use planning.

7         While a number of research studies have been conducted to analyze mobile or location  
8 based social network data, the research is still in its infancy. The analysis has been focused either  
9 on disaggregate patterns or aggregate analysis of social and temporal attributes. While these efforts  
10 provide useful insights, linkages to transportation planning outcomes such as socio-demographics,  
11 land use variables, transportation and infrastructure variables, and points of interests are poorly  
12 understood. The main objective of our proposed effort is to employ check-in data aggregated at  
13 the census tract level to develop a quantitative model for activity intensity as a function of socio-  
14 demographics, transportation infrastructure, land use and built environment attributes. The study  
15 also recognizes that developing a single model for NYC would be restrictive and of limited use.  
16 Hence, prior to modeling, we classify the census tracts in NYC into four clusters as a function of  
17 land use variables. Subsequently, for each cluster a Negative Binomial Regression model is  
18 developed to study activity intensity across the city. The results from these models are employed  
19 to conduct a hot spot analysis highlighting the impact of independent variables across the urban  
20 region. The hot spot analysis illustrated how the data from LBSN users can assist planners to make  
21 informed decisions on mobility and infrastructure needs.



### 3. DATA SOURCE AND DESCRIPTIVE STATISTICS

The original check-in dataset used in a previous study by Cheng et al. (9) was employed in this research. The data consisted of 220,000 unique users checked-in at 1,200 venues from December 2011 to April 2012<sup>1</sup>. The data was obtained from Location Sharing Services (LSS) applications such as Foursquare, Twitter, TweetDeck, Gowalla. Up to 2,000 most recent geo-labeled tweets for each user were saved (for more details on the dataset format see (9)). Using GIS analysis procedures, the check-in data in the NYC region were selected. NYC's population in 2011 was 8.273 million with 2,166 census tract zones based on the zoning system of the US Census Bureau. The aggregated check-in counts were augmented with census tract characteristics including socio-demographics, land-use characteristics, and points of interests. After the data processing, 624,595 geo-coded check-ins were considered for analysis. The check-ins in the census tract range from 0 through 11,159 with an average of about 288.

In our analysis, we generated a host of variables from four broad categories including: (1) land use characteristics (such as one and two family buildings, multi-family walk-up buildings, multi-family elevator buildings, residential buildings, commercial and office buildings, industrial and manufacturing, transportation and utility, public facilities and institutions, open space and outdoor recreation, and parking facilities), (2) socio-demographics (such as population density, population by age, gender, race, and household characteristics), (3) points of interest (locations such as leisure, tourism, recreational, library, airport, sidewalk café, health places), and (4) built environment (such as bus line, bus stops, subway stops, train stops, ferry landing, park and ride stations, bike route, street center line, school counts, building footprint, and green area). We have

---

<sup>1</sup> The proposed prediction framework of activity check-ins can be employed by using LBSN data from other regions or for any other year if the data is available.

1 extracted the abovementioned variables at the Census Tract level for the year 2010 to reflect the  
 2 available LBSN data. A descriptive summary of the characteristics generated for our analysis are  
 3 presented in Table 1.

4  
 5 **TABLE 1 Descriptive Statistics of NYC Census Tracts**

Variables Name	Definition	Zonal		
		Minimum	Maximum	Average
<b>Dependent variable</b>				
Check-in Counts per CT*	Total number of check-ins per CT	0	11159	288.41
<b>Socio-Demographic Characteristics</b>				
Median Age	Median CT Age / 10	0	8.45	3.57
Caucasian Proportion	Caucasian Population of CT / Total Population of CT	0	1.00	0.43
African – American Proportion	African-American population of CT / Total population of CT	0	1.00	0.27
Hispanic Proportion	Hispanic population of CT / Total population of CT	0	1.00	0.26
Asian Proportion	Asian population of CT / Total population of CT	0	1.00	0.12
Children in HH <sup>x</sup>	Total number of children of CT / Total number of HH of CT	0	0.64	0.29
Family HH	Total number of family HH of CT / Total number HH of CT	0	1.00	0.64
Average Family Size	Average family size of CT	0	6.09	3.27
Rental Vacancy Rate (%)	Rental vacancy units*100 / Total number of units within CT	0	61.20	4.75
CT Area	CT area in acres	0.0134	4502.98	89.34
Total Population	Total population per CT	0	26588	3778.70
<b>Points of Interest Characteristics</b>				
Automotive	Number of Automotive Related Places per CT	0	7	0.04
Government	Number of Governmental Places per CT	0	75	4.50
Leisure	Number of Leisure Points per CT	0	13	0.72
Tourism	Number of Touristic Places per CT	0	16	0.07
Health	Number of Health Related Places per CT	0	8	0.12
Library	Number of Libraries per CT	0	2	0.10
Nursing	Number of Nursing Places per CT / 10	0	12	0.56
Senior	Number of Senior places per CT	0	2	0.12
Airport	Number of Airports per CT	0	1	0.00
Ferry Landing	Number of Ferry Landing per CT	0	4	0.02
Beach, Garden, Natural state parks	Number of beaches, gardens, natural state parks per CT / 10 <sup>2</sup>	0	422	2.59

Subway Yard	Number of subway yards per CT	0	2	0.02
Food Places	Number of food aid related places per CT	0	10	0.52
Other Transportation Fac.	Number of other transportation facilities per CT / 10	0	17	0.09
Waste Management Fac.	Number of waste management facilities per CT / 10	0	47	0.24
Children Daycare	Number of children daycare points per CT / 10	0	11	1.19
Bus Depot	Number of bus depots per CT	0	1	0.01
Recreation, Plaza, Mall	Number of recreation areas, plazas, malls per CT / 10	0	12	0.51
Sidewalk Café	Number of sidewalk café per CT / 10 <sup>3</sup>	0	13.60	0.93
<b>Transportation Infrastructure Characteristics</b>				
Street Centerline	Total length of street center lines in ft in CT / 10 <sup>4</sup>	0	90.01	5.38
Bus line	Total length of bus lines in ft in CT / 10 <sup>3</sup>	0	14.47	1.07
Building Footprint	Total area of building footprints in CT / 10 <sup>7</sup>	0	0.87	0.29
Building Elevation	Total number of building floors in CT / 10 <sup>3</sup>	0	6.08	1.19
Green Area Density	Total green area by CT area in sq-ft / 10 <sup>3</sup>	0	6597.24	3.60
Railroad	Total length of railroads in ft in CT / 10 <sup>5</sup>	0	22.92	0.36
Bike Route	Total length of bike routes in ft in CT / 10 <sup>4</sup>	0	6.09	0.20
Number of Buildings	Total number of buildings in CT / 10 <sup>3</sup>	0	3.25	0.49
Bus Stop	Total number of bus stops in CT / 10	0	6	0.61
Subway Entrances	Total number of subway entrances in CT / 10	0	3.60	0.08
Subway Stops	Total number of subway stops in CT	0	6	0.22
<b>Land Use Characteristics</b>				
One and Two Family Buildings Density	Total one and two family building lands by CT area in sq-ft / 10	0	6.79	1.78
Multi-Family Walk-Up Buildings Density	Total multi-family walk-up buildings lands by CT area in sq-ft / 10	0	5.42	0.96
Multi-Family Elevator Buildings Density	Total multi-family elevator buildings lands by CT area in sq-ft / 10	0	11.86	0.72
Residential and Comm. Buildings Density	Total residential and commercial lands by CT area in sq-ft / 10	0	8.64	0.51
Commercial and Office Buildings Density	Total commercial and office buildings lands by CT area in sq-ft / 10	0	6.12	0.39
Industrial and Manufacturing Density	Total industrial and manufacturing lands by CT area in sq-ft / 10	0	5.90	0.19
Transportation and Utility Density	Total transportation and utility lands by CT area in sq-ft / 10	0	8.92	0.16
Public Facilities and Institutions Density	Total public facilities and institution lands by CT area in sq-ft / 10	0	9.10	0.57
Open Space and Outdoor Recreation Density	Total open space and outdoor recreation lands by CT area in sq-ft / 10	0	27.37	0.39
Parking Facilities Density	Total parking facility lands by CT area in sq-ft / 10	0	1.58	0.10

1 \*Census Tract

2 ×Household

## 4. EMPIRICAL ANALYSIS

### (a) 4.1 Cluster Analysis

Clustering is a widely used statistical analysis tool to categorize items together based on their similarities or dissimilarities (18). The aim of clustering algorithm is to classify the population into k categories based on a multivariate set of exogenous variables. The clusters generated should be internally homogenous while being heterogeneous relative to other clusters (19). k-means clustering is a common and straightforward model that uses minimum Euclidean distance between observations (see (20) for similar examples).

Based on the host of exogenous variables, a cluster analysis is conducted to categorize the census tracts. Specifically, eight land use characteristics were employed to undertake the clustering exercise. These characteristics are; one and two-family buildings, multi-family walk-up buildings, multi-family elevator buildings, mixed residential and commercial buildings, commercial and office buildings, industrial and manufacturing, transportation and utility, public facilities and institutions. The k-means clustering algorithm provided good fit for a 4-cluster classification. Also, Bonferroni post-hoc test was used to validate the results of k-means clustering and multiple pairwise comparisons obtained with over 75% of each cluster was found to be statistically significant. The characteristics of the final clusters obtained are presented in Table 2. The spatial distribution of census tracts identified with cluster analysis results are illustrated in Figure 1. Cluster 1 consists of census tracts with single family and/or townhouses. Cluster 2 is represented by a mix of public facilities, commercial buildings and offices, transportation as well as elevated/high rise buildings. Cluster 3 is composed of low rise residential and commercial buildings. Finally, as can be seen from the figure, a large portion of Cluster 4 are census tracts surrounding central park in Manhattan area. This cluster is predominantly covered by high rise

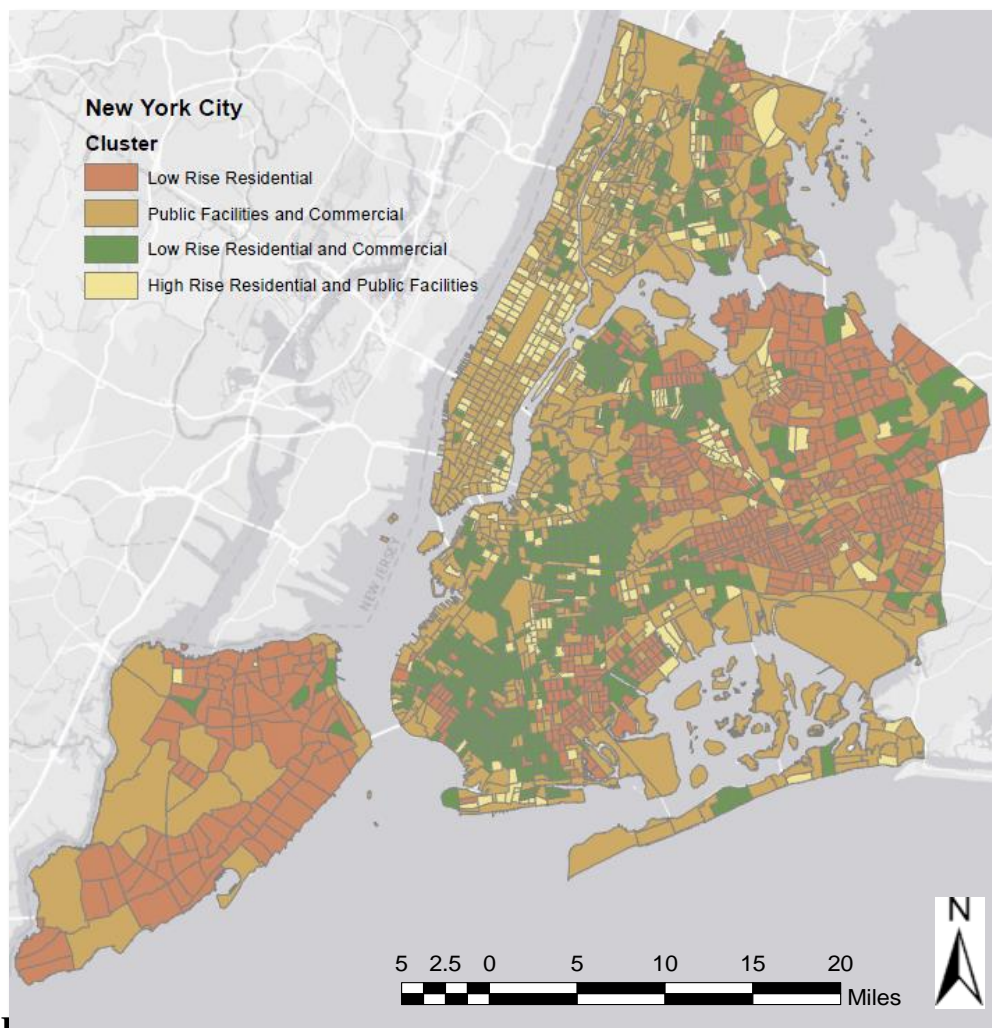
1 buildings in addition to public facilities and institutions. Based on the characteristics, the clusters  
 2 are labelled as follows: Cluster 1 – Low-rise residential, Cluster 2 – Public facilities and  
 3 Commercial, Cluster 3 – Low-rise residential and commercial, Cluster 4 – High-rise residential  
 4 and public facilities (see Figure 1).

5

6 **TABLE 2 Final Cluster Centers**

<b>Land Use Variables</b>	<b>Cluster 1 (Low Rise Residential)</b>	<b>Cluster 2 (Public Facilities and Commercial)</b>	<b>Cluster 3 (Low Rise Residential and Commercial)</b>	<b>Cluster 4 (High Rise Residential and Public Facilities)</b>
One and Two Family Buildings	4.01	0.50	1.65	0.43
Multi-Family Walk-Up Buildings	0.47	0.54	1.95	0.68
Multi-Family Elevator Buildings	0.13	0.56	0.38	3.30
Mixed Residential and Commercial Buildings	0.17	0.72	0.53	0.72
Commercial and Office Buildings	0.21	0.70	0.29	0.31
Industrial and Manufacturing	0.08	0.44	0.11	0.06
Transportation and Utility	0.06	0.37	0.09	0.08
Public Facilities and Institutions	0.30	0.90	0.50	0.57
Number of Census Tract Zones	587	664	650	265

7



1  
2 **FIGURE 1 Spatial Distribution of Clusters**

3  
4 **(b) 4.2 Negative Binomial (NB) Regression Model Results**

5 Given that activity intensity is represented based on non-negative integers, Negative Binomial  
6 (NB) regression approach is employed for our analysis. For the sake of brevity, details on the  
7 model formulation are not provided (see (21) for more details). For model estimation, two sets of  
8 models were estimated. First, a single NB model for New York City CT's (census tracts) was  
9 developed (pooled model). Second, NB models specific to each cluster (obtained above) were  
10 estimated.

1           Prior to discussing the estimation results, we compare the performance of the pooled model  
2 and cluster models. The model performance was tested based on the computation of Bayesian  
3 Information Criterion (BIC) that penalizes the model with large number of parameters. The BIC  
4 for a given empirical model is equal to:  $BIC = -2LL + K \ln(Q)$ ; where  $LL$  is the log likelihood  
5 value at convergence,  $K$  is the number of parameters, and  $Q$  is the number of observations. The  
6 model with the lower BIC is the preferred model. The corresponding BIC values for pooled and  
7 cluster models are: 23376.2 and 23304.3, respectively. The comparison clearly illustrates the  
8 improved fit offered by the cluster specific NB models. For the sake of brevity, we restrict  
9 ourselves to the discussion of cluster-based NB models. The model estimation results for the  
10 cluster-based NB models are presented in Table 3.

11

### 12 (c) *Socio-Demographic Characteristics*

13 Several sociodemographic characteristics influence the activity intensity at the census tract level  
14 including: median age by gender, proportion of population by ethnicity, average number of  
15 children at the household level, average family size, the proportion of family households within  
16 census tracts and percentage of rental vacancy rate.

17           Across the four clusters, the increase in median age has an overall negative effect. While  
18 median age coefficients by gender are positive (male median age for cluster 2 and 3 or female  
19 median age for cluster 4), the other median age coefficient is negative and slightly larger in  
20 magnitude. The result confirms the finding of Sloan et al. (22) that increasing median age in the  
21 census tract reduces activity intensity. Proportion of population by ethnicity has varying trends  
22 across clusters. In cluster 1, higher proportion of Caucasian and African-American increases  
23 activity intensity. On the other hand, for cluster 2, higher proportion of African-American and

1 Hispanic ethnicities are likely to reduce check-in activity. In cluster 3, Hispanic proportion has a  
2 positive influence while Caucasian proportion has a positive influence in cluster 4. On the other  
3 hand, according to Pew Research Center (*1*), overall statistics indicate that Hispanic and African-  
4 American social media users proportion are slightly higher than Caucasian users proportion. The  
5 findings provide evidence that the same variable can affect census tracts across the region  
6 differently. These trends could not have been captured using a pooled model. The increase in the  
7 average number of children at the household level, as expected, reduces the activity intensity across  
8 clusters with varying magnitudes. The presence of higher proportion of family households reduces  
9 activity in Cluster 3. Correspondingly, a study effort by Do et al. (*23*) observed that family  
10 households have lower average weekly visited places compared to other household types. The  
11 increase in average family size has a positive influence on activity intensity for Cluster 2. Finally,  
12 the rental vacancy rate has a negative influence on Cluster 3, probably because the increase in  
13 rental vacancy represents lower occupancy rate resulting in lower number of activities.



1 **Table 3 Negative Binomial Regression Results**

Variable Names*	Cluster 1		Cluster 2		Cluster 3		Cluster 4	
	Estimate	t-stat	Estimate	t-stat	Estimate	t-stat	Estimate	t-stat
Constant	6.32	5.65	6.12	21.17	8.20	13.10	8.97	12.13
<b>Socio-Demographic Characteristics</b>								
Median Age	-0.67	-3.77	-0.27	-5.78	-0.38	-3.13	-0.74	-5.79
Caucasian Proportion	2.26	4.50	-	-	-	-	0.790	2.21
African – American Proportion	1.48	2.98	-0.61	-3.99	-	-	-	-
Hispanic Proportion	-	-	-0.73	-3.97	0.55	3.04	-	-
Asian Proportion	2.65	4.34	-	-	-	-	-	-
Children in HH	-12.91	-7.62	-8.36	-13.01	-5.90	-5.41	-9.30	-8.15
Family HH	-	-	-	-	-4.01	-6.29	-	-
Avg. Family Size	-	-	0.39	5.94	-	-	-	-
Rental Vacancy Rate (%)	-	-	-	-	-0.05	-3.63	-	-
<b>Points of Interest Characteristics</b>								
Automotive	-	-	-	-	-	-	1.22	2.68
Government	-	-	0.03	5.43	-	-	-	-
Leisure	-	-	0.06	2.23	-	-	0.24	3.71
Tourism	-	-	-	-	-	-	0.87	3.20
Health	-	-	-	-	-0.25	-2.91	-	-
Library	-	-	-0.27	-2.49	-	-	-	-
Nursing	-	-	-	-	0.12	2.48	-	-
Senior	0.75	3.43	-	-	-	-	-	-
Airport	-	-	0.31	4.49	-	-	-	-
Ferry Landing	-	-	-	-	-	-	0.92	2.20
State Parks, National and Cultural Inst.	-	-	-	-	-	-	-0.43	-1.99
Food Places	0.47	4.92	-	-	-0.09	-2.64	-	-
Other Transportation Fac.	-	-	-	-	0.32	2.52	-1.55	-2.19
Waste Management Fac.	-	-	-	-	-0.29	-3.35	0.35	2.72
Children Daycare	-	-	-	-	0.04	1.69	0.08	2.03
Bus Depot	-	-	-	-	-	-	-2.84	-2.23
Recreation, Plaza, Mall	0.16	2.56	0.05	2.08	-	-	-	-
Sidewalk Café	0.18	1.82	0.08	2.00	-	-	-	-
<b>Transportation Infrastructure Characteristics</b>								
Street Centerline	0.07	3.48	0.03	4.55	0.08	2.80	-	-
Bus line	0.08	2.14	-	-	-	-	-	-
Railroad	-	-	0.05	2.24	-	-	0.18	1.74
Bike Route	-	-	0.39	5.65	-	-	-	-

Bus Stop	-	-	0.12	1.86	0.18	1.69	-	-
Subway Entrances	-0.99	-1.76	-	-	1.45	4.72	1.44	3.21
<b>Land Use Characteristics</b>								
One and Two Family Buildings Density	0.12	2.58	-0.61	-7.84	-	-	-0.61	-6.28
Multi-Family Walk-Up Build. Density	0.242	2.05	-	-	-	-	-0.17	-1.81
Residential and Comm. Build. Density	-	-	0.129	2.74	0.516	4.52	-	-
Commercial and Office Build. Density	1.06	5.38	0.28	5.88	0.43	3.78	-	-
Industrial and Manufacturing Density	0.36	1.74	-0.13	-2.60	-	-	-	-
Transportation and Utility Density	-	-	-	-	0.67	3.73	-	-
Building Footprint	-	-	1.35	3.05	-	-	-	-
Building Elevation	0.18	2.29	0.17	2.35	1.31	4.55	0.93	6.06
Green Area Density	-	-	-	-	-	-	-0.06	-3.06
Number of Buildings	-	-	-	-	-1.70	-2.81	-	-
Parking Facilities Density	-	-	0.36	2.25	-0.06	-1.78	-	-
<b>Summary Statistics</b>								
Number of Census Tracts	587		664		650		265	
Log-Likelihood	-2551.57		-4080.90		-3245.30		-1486.77	
LR chi square (Number of Predictors)	395.12 (17)		1151.98 (22)		793.80 (20)		323.38 (18)	
Pseudo R <sup>2</sup>	0.072		0.12		0.11		0.09	

1 \* Variable definitions are presented in Table 1

1     (d)     *Points of Interest Characteristics*

2     Several points of interest characteristics influence the activity intensity observed. Cluster 1 is  
3     positively influenced by senior centers, food places, recreation plaza and malls, and sidewalk cafes.  
4     According to the venue cloud for check-ins generated by Cheng et al. (9), it is clear that the largest  
5     clouds are café's (i.e. coffee shop), food places, and centers (i.e. shopping malls). In Cluster 2,  
6     government related, leisure related, airport recreation plaza and mall, and sidewalk cafes have a  
7     positive influence while libraries have a negative influence. Li et al. (24) indicated that a large  
8     proportion of users are likely to check-in at particular places such as airports. For cluster 3, the  
9     variable affecting activity intensity positively include nursing related, other transportation  
10    facilities, children day care and sidewalk facilities. On the other hand, variables affecting  
11    negatively include health related, food places, and waste management facilities. Finally, for  
12    Cluster 4, automotive related, leisure related, tourism related, ferry landing, beach, garden and  
13    natural facilities, waste management, children day care show positive influence. Other  
14    transportation facilities, and bus depot affect activity negatively in Cluster 4. Overall, the results  
15    capture the variation across the various clusters based on the points of interest. The results are hard  
16    to compare to earlier work because detailed information of this resolution has rarely been  
17    employed in transportation planning applications.

18

19    (e)     *Transportation Infrastructure*

20    The impact of transportation infrastructure offers significant differences across the clusters. The  
21    street centerline length has a positive association with activity intensity in clusters 1 through 3.  
22    The bus line length in the census tract has a positive effect on cluster 1 activity intensity. In  
23    contrast, Sengstock et al. (25) and Frias-Martinez et al. (26) imply that national parks are highly  
24    associated with social media check-ins. The length of rail road has a positive impact on activity

1 intensity for cluster 2 and 4. The bicycle route length variable affects positively the intensity in  
2 cluster 2 only. The number of building variable in cluster 3 has a negative impact on activity  
3 intensity. The number of bus stops has a positive influence on cluster 3 ridership. Finally, number  
4 of subway entrances has a negative influence on cluster 1 activity intensity while positively  
5 influencing activity intensity in clusters 3 and 4. According to the tweet content models developed  
6 for NYC, Kling et al. (13) indicated that transportation facilities are highly mentioned in the  
7 morning period while in Manhattan and East Village they were highly references in the evening  
8 period.

9

#### 10 (f) *Land Use Characteristics*

11 Land use characteristics in the census tracts exhibit significant influence on activity intensity. A  
12 higher density of one and two-family buildings has a positive effect on activity intensity in cluster  
13 1 while reducing activity intensity in clusters 2 and 4. The increase in density of multi-family  
14 walkup units has a positive effect on cluster 1 activity intensity. The residential and commercial  
15 density variable has a positive influence on activity intensity for cluster 2 and 3. Commercial and  
16 office building density is associated with positive influence for clusters 1 through 3. Similarly, Hu  
17 et al. (27) indicates that commercial zones attract people's attention. Industrial and manufacturing  
18 density has a positive influence on activity intensity for cluster 1 and a negative influence on  
19 activity intensity for cluster 2. This finding might be affected by the fact highlighted by Frias-  
20 Martinez et al. (26), that industrial land use is at a minimum in most regions of NYC (i.e. less than  
21 8% in Manhattan). Transportation and utility density are positively associated with cluster 3  
22 activity intensity. Building footprint significantly increases activity intensity for clusters 2.  
23 Building elevation increase is associated with higher activity intensity for all clusters (except 2).

1 Interestingly, green area density is negatively associated with activity intensity in cluster 4. The  
2 building density of the area affects the type of businesses and accordingly affects the behavior of  
3 people visiting these areas (8). Finally, parking facility density has a positive influence on check-  
4 in activity for cluster 2.

5

### 6 (g) 4.3 Model Validation

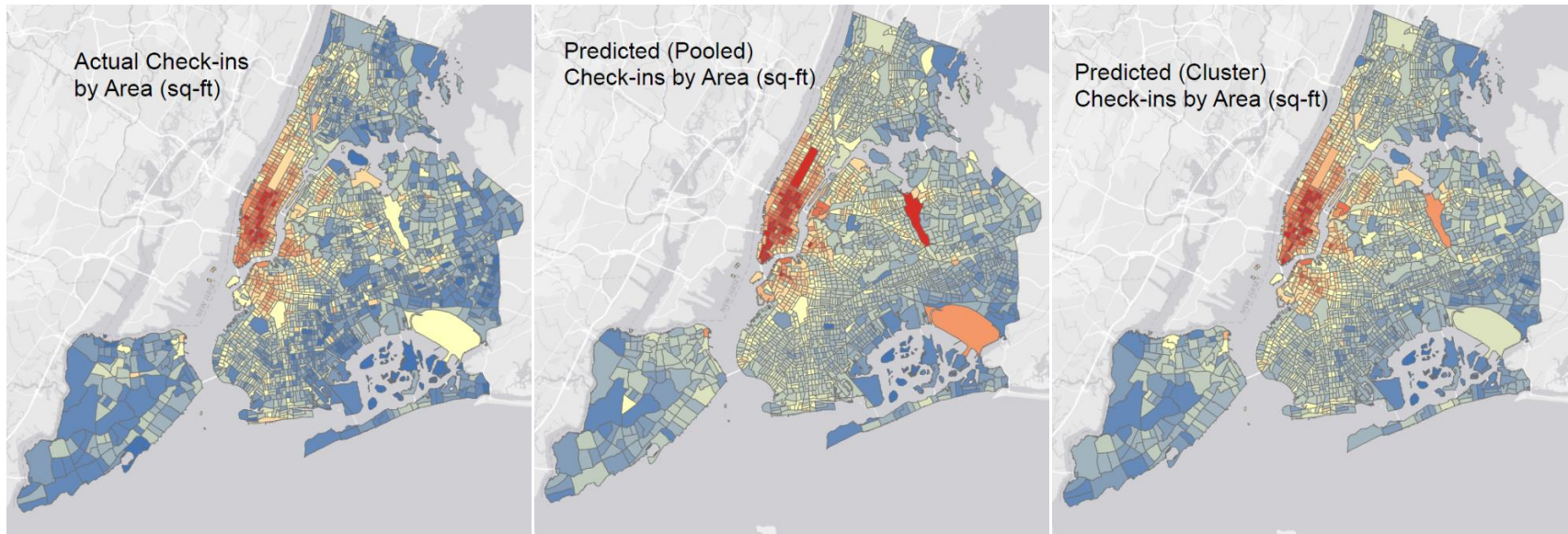
7 To validate the model performance, we spatially represent (a) observed check-ins per unit area, (b)  
8 check-ins per unit area based on pooled model and (c) check-ins per unit area based on cluster-  
9 based models. The patterns of activity check-ins are presented in Figure 2. The categories  
10 considered for the three figures are: 0 - 0.05, 0.05 – 0.3, 0.3 – 0.5, 0.5 – 1, 1 – 3, 3 – 5, 5 – 10, 10  
11 – 25, 25 – 50, 50 – 100, 100 – 200, 200 – 250, 250 and higher. From the visual comparison, across  
12 the three patterns, it is evident that the activity check-in patterns for cluster-based models are closer  
13 to the observed patterns. For instance, the pooled model over-predicts activity around central park  
14 and John F. Kennedy airport while the cluster-based models are closer to the observed patterns.  
15 To be sure, the cluster-based model also produces slightly different estimates for some census  
16 tracts. But overall, it offers more close resemblance to observed patterns.

17

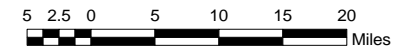
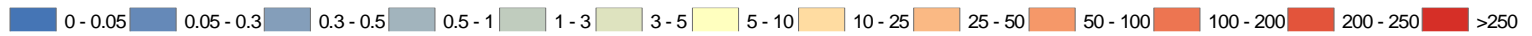
### 18 (h) 4.4 Hot Spot Analysis

19 In this section, to illustrate the influence of exogenous variables, we undertake a unique hot spot  
20 analysis. The hot spot analysis is based on the value of the contribution of the individual parameter  
21 to the count propensity ( $\beta * x_n$ ). The contribution to count propensity is plotted by implementing  
22 Optimized Kernel Density tool of GIS ArcMap. The tool automatically aggregates the predicted  
23 check-in frequency, identifies an appropriate scale of analysis, and corrects for both multiple

1 testing and spatial dependence by calculating the mean center of the input points using a radius  
2 search (bandwidth) algorithm. This tool allows us to identify statistically significant spatial groups  
3 of high values (hot spots) and low values (cold spots). Statistically significant hot and cold spots  
4 indicate that rather than a random pattern, the corresponding explanatory variable prediction  
5 exhibit statistically significant spatial dispersion. The variables chosen for the hot spot analysis  
6 are: Children by HH, One and Two-Family Buildings, Sidewalk Café, Median Age, Street  
7 Centerline and Building Elevation. The spatial representations are presented in Figure 3. Light  
8 green background color indicates that both hot and cold spots exist on the figure, whereas blue  
9 background indicates that the heat map includes only hot or cold spots along with neutral areas.  
10 The results clearly illustrate the spatial regions that are significantly affected by these variables  
11 across the NYC region.

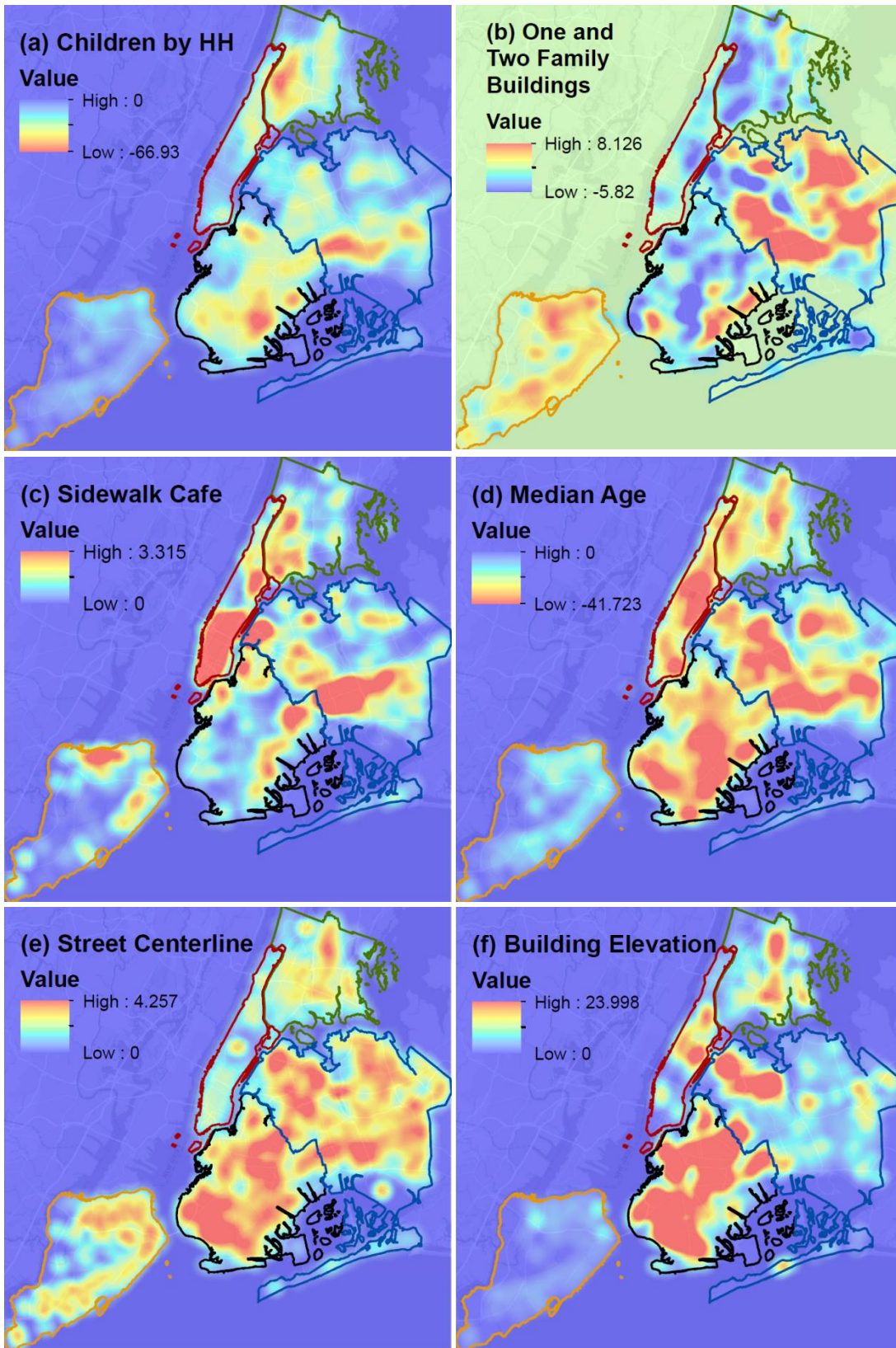


*Check-in Density Categories:*



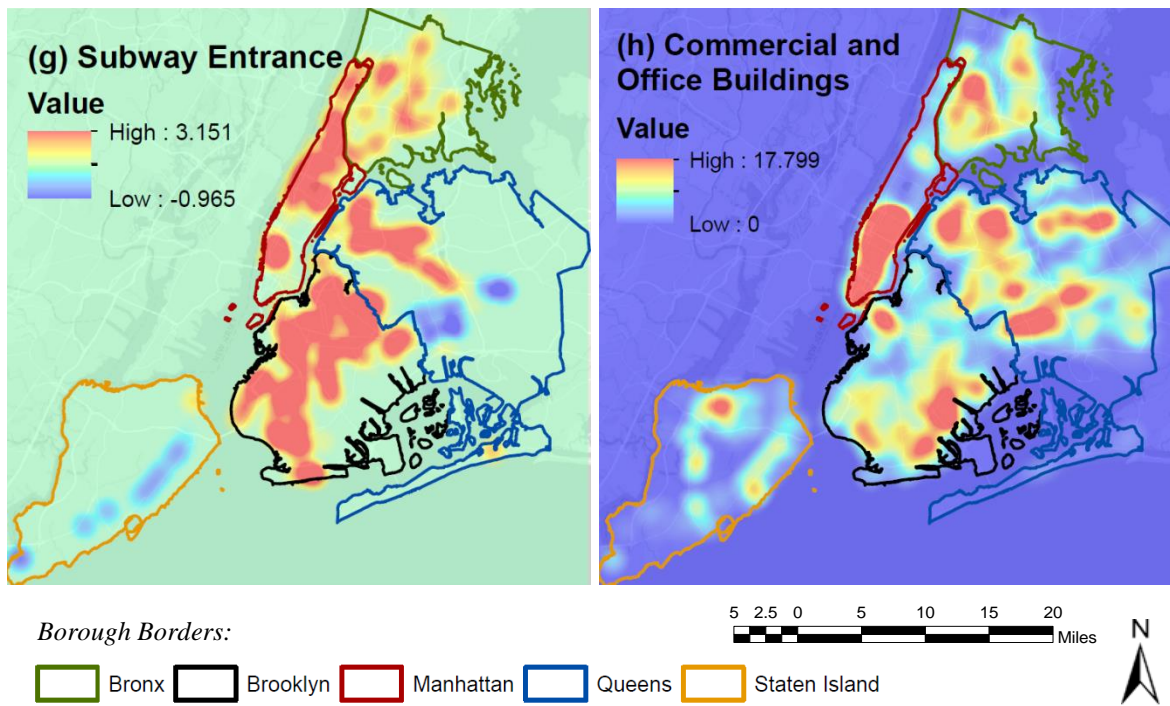
1  
2

**FIGURE 2 New York City Check-in Density Predictions**



1  
2





1  
 2 **FIGURE 3** Kernel Density Estimate Heat Maps for Selected Exogenous Variables; (a) Children  
 3 by HH, (b) One and Two Family Buildings, (c) Sidewalk Café, (d) Median Age, (e) Street  
 4 Centerline, (f) Building Elevation, (g) Subway Entrance, (h) Commercial and Office Buildings  
 5

## 6 5. CONCLUSION

7 The current study employed a location-based social networks (LBSN) service-based data for  
 8 aggregate level transportation planning exercise by developing land-use planning models.  
 9 Specifically, we employed check-in data aggregated at the census tract level to develop a  
 10 quantitative model for activity intensity as a function of land use and built environments attributes  
 11 for the New York City (NYC) region. The detailed exogenous variables considered were socio-  
 12 demographics, land use variables, transportation variables, and points of interests at the census  
 13 tract level. The study also recognized that developing a single model for NYC would be restrictive  
 14 and of limited use. Hence, prior to modeling, we classified the census tracts in NYC into four  
 15 groups as a function of eight different land use variables. The clusters identified were labelled as  
 16 follows: Cluster 1 – Low-rise residential, Cluster 2 – Public facilities and Commercial, Cluster 3

1 – Low-rise residential and commercial, and Cluster 4 – High-rise residential and public facilities.  
2 The clustering approach, rather than considering the entire city as homogenous allowed us to  
3 distinguish across different clusters. Subsequently, for each cluster as well as for the whole region,  
4 Negative Binomial (NB) Regression models were developed to study activity intensity patterns  
5 across the city. We compared the performance of the pooled model and cluster models by using  
6 Bayesian Information Criterion. The comparison clearly illustrated the improved fit offered by the  
7 cluster specific NB models.

8         From the estimation results, we found that there are variations across the different clusters  
9 based on different exogenous variables. Moreover, the variables effects found to be different for  
10 some clusters and the pooled model supporting our hypothesis that activity intensity profile is not  
11 same across the entire region. To further validate the model performance, we spatially represented  
12 the observed check-ins and predicted check-ins based on pooled and cluster-based models. From  
13 the visual comparison, across the three patterns, it was evident that the activity check-ins pattern  
14 for cluster-based models is closer to the observed patterns. We also illustrated the impact of various  
15 parameters on check-ins using a hot spot analysis. This tool enabled us to identify statistically  
16 significant spatial groups of high values (hot spots) and low values (cold spots). The results clearly  
17 illustrated the spatial regions that are significantly affected by different variables across the NYC  
18 region. The findings from our study provided insights on relative differences of activity  
19 engagements across the urban region. The proposed approach thus provides a complementary  
20 analysis tool to traditional transportation planning exercises.

21         The paper is not without limitations. The dataset employed in our analysis is from  
22 December 2011 through April 2012. Ideally, the consideration of a more recent time would be  
23 beneficial. The reader would note that the methodology developed could be applied to analyze

1 newer versions of data that are freely available or purchased at a cost for an urban region of interest.  
2 The data used in our analysis is for a part of the year. Hence, accommodating for seasonality effects  
3 was not possible. Towards accommodating for these effects, it would be useful to consider  
4 obtaining data for a full year and generating the Check-in measures across different seasons. The  
5 dependent variable thus generated can be analyzed using the proposed model to identify  
6 seasonality differences. For our analysis, we did not consider the spatial correlations across  
7 different neighboring census tracts. In the future, it might be beneficial to examine for the influence  
8 of spatial correlation in the count models.

## 9 (i) REFERENCES

- 10 1. P. R. Center, "Social media fact sheet," *Pew Res. Cent. Internet, Sci. Tech*, 2017.
- 11 2. E. Gordon and A. de S. e Silva, *Net locality: Why location matters in a networked world*.  
12 John Wiley & Sons, 2011.
- 13 3. S. Hasan and S. V. Ukkusuri, "Urban activity pattern classification using topic models from  
14 online geo-location data," *Transportation Research Part C Emerging Technologies*, vol.  
15 44, pp. 363–381, 2014.
- 16 4. M. Rzeszewski, "Geosocial capta in geographical research—a critical analysis,"  
17 *Cartography and Geographic Information Science*, vol. 45, no. 1, pp. 18–30, 2018.
- 18 5. V. Frias-Martinez, V. Soto, H. Hohwald, and E. Frias-Martinez, "Characterizing urban  
19 landscapes using geolocated tweets," *Proceedings - 2012 ASE/IEEE International  
20 Conference on Privacy, Security, Risk and Trust and 2012 ASE/IEEE International  
21 Conference on Social Computing, SocialCom/PASSAT 2012*, 2012, pp. 239–248.
- 22 6. A. R. Pinjari and C. R. Bhat, "Activity-based travel demand analysis," in *A Handbook of  
23 Transport Economics*, Edward Elgar Publishing, 2011.

- 1 7. H. J. Miller, "Activity-based analysis," *Handbook of regional science*, Springer, 2014, pp.  
2 705–724.
- 3 8. J. Cranshaw, J. I. Hong, and N. Sadeh, "The Livelihoods Project: Utilizing Social Media to  
4 Understand the Dynamics of a City," *The 6th International AAAI Conference on Weblogs  
5 and Social Media*, pp. 58–65, 2012.
- 6 9. Z. Cheng, J. Caverlee, K. Lee, and D. Z. Sui, "Exploring Millions of Footprints in Location  
7 Sharing Services," *Icwsn*, vol. 2010, no. Cholera, pp. 81–88, 2011.
- 8 10. R. Ahas, A. Aasa, Y. Yuan, M. Raubal, Z. Smoreda, Y. Liu, C. Ziemlicki, M. Tiru, and M.  
9 Zook, "Everyday space–time geographies: using mobile phone-based sensor data to  
10 monitor urban activity in Harbin, Paris, and Tallinn". *International Journal of  
11 Geographical Information Science*, vol. 29, no 11, pp.2017-2039, 2015.
- 12 11. G. Cao, S. Wang, M. Hwang, A. Padmanabhan, Z. Zhang, and K. Soltani, "A scalable  
13 framework for spatiotemporal analysis of location-based social media data," *Computers,  
14 Environment and Urban Systems*, vol. 51, pp. 70–82, 2015.
- 15 12. A. Ahmed, L. Hong, and A. J. Smola, "Hierarchical geographical modeling of user  
16 locations from social media posts," *Proceedings of the 22nd international conference on  
17 World Wide Web*, 2013, pp. 25–36.
- 18 13. F. Kling and A. Pozdnoukhov, "When a City Tells a Story: Urban Topic Analysis,"  
19 *Proceedings of the 20th International Conference on Advances in Geographic Information  
20 Systems - SIGSPATIAL '12*, p. 482, 2012.
- 21 14. S. Wakamiya, R. Lee, and K. Sumiya, "Urban area characterization based on semantics of  
22 crowd activities in Twitter," (*Including Subseries Lecture Notes in Artificial Intelligence  
23 and Lecture Notes in Bioinformatics*), 6631 LNCS, 108–123. <https://doi.org/10.1007/978->

- 1           3-642-20630-6\_7ecture Notes in Computer Science, vol. 6631 LNCS, pp. 108–123, 2011.
- 2   15.   A. Noulas, S. Scellato, R. Lambiotte, M. Pontil, and C. Mascolo, “A tale of many cities:  
3           Universal patterns in human urban mobility,” *PLoS One*, vol. 7, no. 5, 2012.
- 4   16.   A. Bawa-Cavia, “Sensing the urban: using location-based social network data in urban  
5           analysis,” *Pervasive PURBA Workshop*, 2011.
- 6   17.   X. Zhan, S. V. Ukkusuri, and F. Zhu, “Inferring Urban Land Use Using Large-Scale Social  
7           Media Check-in Data,” *Networks Spat. Econ.*, vol. 14, no. 3–4, pp. 647–667, 2014.
- 8   18.   M. R. Anderberg, *Cluster analysis for applications: probability and mathematical  
9           statistics: a series of monographs and textbooks*, vol. 19. Academic press, 2014.
- 10  19.   M. G. Karlaftis and A. P. Tarko, “Heterogeneity considerations in accident modeling,”  
11           *Accident Analysis & Prevention*, vol. 30, no. 4, pp. 425–433, 1998.
- 12  20.   A. Pinjari, N. Eluru, C. Bhat, R. Pendyala, and E. Spissu, “Joint model of choice of  
13           residential neighborhood and bicycle ownership: accounting for self-selection and  
14           unobserved heterogeneity,” *Transp. Res. Rec. J. Transp. Res. Board*, no. 2082, pp. 17–26,  
15           2008.
- 16  21.   R. Winkelmann, *Econometric analysis of count data*. Springer Science & Business Media,  
17           2008.
- 18  22.   L. Sloan, J. Morgan, P. Burnap, and M. Williams, “Who tweets? Deriving the demographic  
19           characteristics of age, occupation and social class from Twitter user meta-data,” *PLoS One*,  
20           vol. 10, no. 3, p. e0115545, 2015.
- 21  23.   T. M. T. Do and D. Gatica-Perez, “The places of our lives: Visiting patterns and automatic  
22           labeling from longitudinal smartphone data,” *IEEE Transportation Mobility Computations*,  
23           no. 1, p. 1, 2013.

- 1 24. L. Li, M. F. Goodchild, and B. Xu, "Spatial, temporal, and socioeconomic patterns in the  
2 use of Twitter and Flickr," *Cartography and Geographic Information Science*, vol. 40, no.  
3 2, pp. 61–77, 2013.
- 4 25. C. Sengstock and M. Gertz, "Latent geographic feature extraction from social media,"  
5 *Sigspatial*, p. 149, 2012.
- 6 26. V. Frias-Martinez and E. Frias-Martinez, "Spectral clustering for sensing urban land use  
7 using Twitter activity," *Engineering Applications of Artificial Intelligence*, vol. 35, pp.  
8 237–245, 2014.
- 9 27. Y. Hu, S. Gao, K. Janowicz, B. Yu, W. Li, and S. Prasad, "Extracting and understanding  
10 urban areas of interest using geotagged photos," *Computers, Environment and Urban  
11 Systems*, vol. 54, pp. 240–254, 2015.

12