

# SOCIAL CONNECTEDNESS IN URBAN AREAS\*

Michael Bailey<sup>†</sup>   Patrick Farrell<sup>‡</sup>   Theresa Kuchler<sup>§</sup>   Johannes Stroebe<sup>¶</sup>

## Abstract

We use anonymized and aggregated data from Facebook to explore the spatial structure of social networks in the New York metro area. We highlight the importance of transportation infrastructure in shaping urban social networks by showing that travel time and travel costs are substantially stronger predictors of social connectedness between zip codes than geographic distance is. We also document significant heterogeneity in the geographic breadth of social networks across New York zip codes, and show that much of this heterogeneity is explained by the ease of access to public transit, even after controlling for socioeconomic characteristics of the zip codes' residents. When we group zip codes with strong social ties into hypothetical communities using an agglomerative clustering algorithm, we find that geographically non-contiguous locations are grouped into socially connected communities, again highlighting that geographic distance is an imperfect proxy for urban social connectedness. We also explore the social connections between New York zip codes and foreign countries, and highlight how these are related to past migration movements.

**JEL Codes:** R3, R4

**Keywords:** Social Connectedness, Agglomeration Externalities, Transportation Infrastructure

---

\*This version: June 21, 2019. We thank the Center for Global Economy and Business at NYU Stern for generous research support. We also thank Ed Glaeser, Paul Goldsmith-Pinkham, Matt Jackson, Stijn van Nieuwerburgh, Nico Stier, and Maisy Wong, as well as seminar and conference participants at Facebook, NYU, the NYC Real Estate Conference, the 2019 Ohlstadt Workshop, the GEA Conference in Frankfurt, and the NBER Summer Institute, for their helpful comments. We thank Drew Johnston, Sung Lee, and Hongbum Lee for outstanding research assistance. This research was facilitated through a research consulting agreement between the academic authors (Farrell, Kuchler, and Stroebe) and Facebook. Bailey is an employee at Facebook.

<sup>†</sup>Facebook. Email: [mcb Bailey@fb.com](mailto:mcb Bailey@fb.com)

<sup>‡</sup>Princeton University. Email: [pwfarrell@gmail.com](mailto:pwfarrell@gmail.com)

<sup>§</sup>New York University, Stern School of Business. Email: [tkuchler@stern.nyu.edu](mailto:tkuchler@stern.nyu.edu)

<sup>¶</sup>New York University, Stern School of Business, NBER, and CEPR. Email: [johannes.stroebe@nyu.edu](mailto:johannes.stroebe@nyu.edu)

Social networks influence many aspects of our lives, with social ties providing access to a wide range of new ideas and employment opportunities (see Granovetter, 2005; Jackson, 2014; Bramouille, Galeotti, and Rogers, 2016). In the context of urban economics, theories of agglomeration feature the ability to learn from many different people as a key force behind the high productivity of cities (e.g., Jacobs, 1969; Bairoch, 1991; Glaeser, 2011). For example, Glaeser, Kallal, Scheinkman, and Shleifer (1992) describe that “the cramming of individuals, occupations, and industries into close quarters provides an environment in which ideas flow quickly from person to person.” In practice, the strength of these positive agglomeration forces depends on the extent to which individuals living in the same city actually interact with one another, in particular across demographic groups and geographic distances. Indeed, it is likely that cities that facilitate interactions across all inhabitants are best positioned to capitalize on the benefits of agglomeration. However, despite this important role of cities’ social structures in creating agglomeration externalities, data challenges in measuring social networks have severely reduced researchers’ ability to study this social structure at a large scale.

In this paper, we investigate the spatial structure of social networks within the New York metro area. We measure social networks using aggregated and anonymized data from Facebook, a global online social network. By the end of 2017, Facebook had 239 million monthly active users in the U.S. and Canada, and about 2.1 billion such users globally. We observe an anonymized snapshot of all Facebook users with location history enabled as of March 2018. For these users, we observe their locations at the zip code level as well as their connections to other individuals on Facebook. We use these data to explore the local, domestic, and international networks of Facebook users in both New York City (NYC) and the wider New York Combined Statistical Area (New York CSA). The density, diversity, and large population of New York, combined with its varied geography and extensive public transportation infrastructure, present an ideal setting for investigating the factors that influence social network structure in urban settings. Indeed, we believe that our study brings the most comprehensive data to date to measure and explore the social structure of cities.<sup>1</sup> Our empirical approach complements an exciting recent literature that has used cell phone call records to better understand the geography of social connectedness (e.g., Schlöpfer et al., 2014; Herrera-Yague et al., 2015; Büchel and von Ehrlich, 2016). Relative to that literature, the Facebook data capture many more links per individual, allowing us to measure the prevalence and distribution of potentially weak ties that have been shown to be important in the dissemination of information and ideas (Granovetter, 1977).<sup>2</sup> While we are unable to make any conclusive causal inferences on the determinants and effects of the observed social structures, we hope that the novel patterns presented in this paper can help advance our understanding of social connectedness in urban areas.

---

<sup>1</sup>The zip code-level social connectedness data that we compile and use in this project is accessible to researchers and policy makers by emailing a 1-page proposal to [sci\\_data@fb.com](mailto:sci_data@fb.com). See Bailey et al. (2018b) for detailed information on county-level social network data that is also accessible to researchers.

<sup>2</sup>In addition, interactions via phone are often substitutes to in-person interactions. One might therefore worry that researchers’ ability to observe a social link in phone records is systematically related to the frequency of the two individuals interacting in person. The latter should correlate both with geographic distance and the ease of travel via public transport.

In the first part of the paper, we explore the role of public transit infrastructure as a potential determinant of social networks in urban areas. We first discuss a number of case studies that show that the social networks of urban zip codes are distributed along transit routes that connect these zip codes to other parts of the city. To explore the relationship between social connectedness and transportation infrastructure more formally, we calculate the travel times on public transit between each pair of NYC zip codes. We find that social connectedness declines strongly in the travel time between locations. Within NYC, the elasticity of social connectedness to travel time is -1.42, which is about 60% larger in magnitude than the elasticity of social connectedness to distance, which is -0.87. This finding suggests that public transit can help facilitate the maintenance and formation of social links across individuals living in geographically distant parts of the same city. As a result, extensive public transportation infrastructure can increase agglomeration benefits as well as reduce the extent to which residential segregation leads to social segregation. This result is consistent with recent findings that suggest that transportation infrastructure allows individuals to visit restaurants that are farther away, thereby lowering the segregation of consumption patterns (Davis et al., 2017)

In addition to the role played by geographic distance and public transit travel time in forming and maintaining social ties between geographies, we find that zip codes that are more similar along demographic measures such as race, education, and income are more likely to be socially connected. This is consistent with previous studies that have documented that social ties are generally more common between similar individuals and regions, a feature that is often referred to as “homophily” (Lazarsfeld and Merton, 1954; Zipf, 1949; Verbrugge, 1983; Marmaros and Sacerdote, 2006; Bailey et al., 2018a,b). We show that short public transit travel times are more important for connecting zip codes with different incomes than they are for connecting zip codes with similar incomes. This finding highlights that public transit investments do not just facilitate social connections between far-away zip codes in general, but do so particularly across zip codes with different demographics.

We next provide a descriptive analysis of the geographic concentration of social networks. We find substantial heterogeneity in this social network concentration across NYC zip codes. For residents of the median zip code, 29.0% of U.S.-based friends live within 5 miles, but this number ranges from 19.5% to 39.6% between the 5th and the 95th percentiles of the zip code distribution. Similarly, for the median NYC zip code, 22.0% of U.S.-based friends live among the nearest 1 million people, while the 5-95 percentile range is 13.1% to 32.7%. Consistent with the results described above, this geographic concentration of social networks is highly correlated with access to public transportation infrastructure (measured, for example, by the share of a zip code’s population that lives within a quarter mile of a rail transit station). These results hold even after conditioning on zip code demographic and income measures. The ease of transit also explains more of the across-zip code variation in the concentration of social networks than zip code demographics do. Quantitatively, a 15 minute increase in the average travel time to all zip codes is associated with a 3.7 percentage point increase in the share of friends living within the nearest 500k people. The geographic concentration of social net-

works also correlates with socioeconomic outcomes such as income and education levels: the share of friends living within various distances is decreasing in zip code income and increasing in the fraction of population without a high school degree. Although our data do not allow us to make statements about the causal connection between social connectedness and socioeconomic outcomes, our findings are consistent with the urban economics literature that points to social interactions as a primary channel for agglomeration externalities that can improve the economic outcomes for residents.

After exploring the determinants of social connectedness between zip code pairs, we run a hierarchical agglomerative linkage clustering algorithm to construct hypothetical communities of zip codes that maximize within-group social connectedness. We find that although all the communities are contiguous at the CSA level, some communities are non-contiguous when focusing on zip codes within NYC. This finding reinforces the earlier observation that geographic distance might not be as relevant a measure to understand social ties within dense urban areas.

In the final part of the paper, we study the social connectedness of New York zip codes to foreign countries. We find strong heterogeneities in the degree to which different zip codes are connected to different countries. We show that past migration movements are a strong determinant of connections abroad, which is suggestive of immigrants' desire to live in areas near the existing ethnic enclaves or areas with transportation accessible to these communities. Therefore, the clustering of ethnicities in a region plays a key role in explaining the presence of international friendship links.

In terms of measurement, our paper contributes to a recent literature that has used data from online services such as Yelp and Twitter to better understand various elements of social and economic activity within cities (e.g., Davis et al., 2017; Glaeser, Kim, and Luca, 2017). We also build on a literature that has studied the unique properties of urban social networks (see Glaeser and Kahn, 2004; Glaeser, 2011; Kowald et al., 2013; Ioannides, 2013; Herrera-Yague et al., 2015; Ioannides, 2015; Picard and Zenou, 2018). Our novel data allow us to document that public transit infrastructure likely is a crucial determinant of the formation and maintenance of social ties in urban areas, in particular across locations with different demographic makeups. This suggests a mechanism through which transit infrastructure affects social network formation, which in turn can influence economic outcomes. In this sense, our work contributes to an important literature that has shown that transit investments generate immediate economic effects and cause long-term changes to the structure of cities. For instance, Perlman (2016) finds that transportation improvements had significant impact on increases in patenting, especially for counties that were not previously well-connected, and Glaeser (2005) finds that New York has become America's largest city due to its initial dominance as a hub of the transportation system (see also Glaeser and Shapiro, 2001; Glaeser and Gottlieb, 2009; Baum-Snow, 2013; Ioannides, 2013; Brooks and Lutz, 2014; Glaeser and Steinberg, 2016). We hope that the increasing availability of social network data from online social networking services such as Facebook will further boost research efforts that explore the determinants and effects of the social structures of cities.



# 1 Data

We construct our measures of the social connectedness across locations using anonymized administrative data from Facebook, a global online social networking service. Facebook was created in 2004, and, by the end of 2017, had 2.1 billion monthly active users globally and 239 million such users in the U.S. and Canada. An independent survey of Facebook users from 2015 found that more than 68% of the U.S. adult population and 79% of online adults in the U.S. used Facebook (Duggan, Greenwood, and Perrin, 2016). That same survey shows that Facebook usage rates among U.S.-based online adults were relatively constant across income groups, education levels, and race, and among urban, rural, and suburban residents; usage rates were slightly declining in age (from 88% of individuals aged 18 to 29, to 62% of individuals aged 65 and older). Establishing a connection on Facebook requires the consent of both individuals, and there is an upper limit of 5,000 on the number of connections a person can have. As a result, Facebook connections are primarily between real-world acquaintances. Indeed, a second independent survey of Facebook users revealed that only 39% of users reported being Facebook friends with someone they had never met in person (Duggan et al., 2015). In contrast, Facebook users generally reported that they were Facebook friends with real-life friends: 91% said they were Facebook friends with current friends and 87% said they were connected to past friends, such as former classmates. Furthermore, most users reported that they were Facebook friends with their family members: 93% of Facebook users said they were Facebook friends with family members other than parents or children, 45% said they were Facebook friends with their parents, and 43% said they were Facebook friends with their children. Finally, Facebook networks often capture other important social ties: 58% of users said that they were Facebook friends with co-workers and 36% of users reported that they were Facebook friends with their neighbors (Duggan et al., 2015). As a result, networks formed on Facebook more closely resemble real-world social networks than those on other online platforms, such as Twitter, where uni-directional links to non-acquaintances, such as celebrities, are common (see Bailey et al., 2017, 2018a,b, 2019, for additional evidence that friendships observed on Facebook serve as a good proxy for real-world U.S. social connections).

We observe an anonymized snapshot of all active Facebook users from March 2018. We focus on those users who had location history enabled, and who had interacted with Facebook over the 30 days prior to the date of the snapshot. We match those users who reside within the New York Combined Statistical Area (CSA) to their zip code locations. The New York CSA consists of 35 counties across the states of Connecticut, New Jersey, New York, and Pennsylvania. We count as within the New York CSA all Census Bureau Zip Code Tabulation Areas (ZCTAs) that fall at least partly within a county making up the New York CSA. Users within the United States but not within the New York CSA are mapped to their county of residence. Users outside of the United States are mapped to their country of residence. From these data, we obtain a count of the number of connections between each zip code  $i$  in the New York CSA and each other region  $j$ , where  $j$  is either another zip code within the New York CSA, a U.S. county outside of the New York CSA, or a foreign country.

We only include zip codes in our analysis that have a total population of at least 500 people and that are above the 5th percentile in the number of eligible Facebook users within the New York CSA. These restrictions are intended to preserve user anonymity as well as to reduce the improper matching of users to officially unpopulated or unusual zip codes, such as individual non-residential buildings (e.g., post offices) or abnormal locations (e.g., JFK airport). Our final data set includes 182 zip codes in NYC and 1,181 zip codes across the entire New York CSA.

We combine these data on social networks with information on the population and demographics of zip codes from the 2015 Census Bureau 5-year American Community Survey (ACS) and the 2014 Internal Revenue Service (IRS) Individual Income Tax Statistics. In particular, information on total population, racial composition, and educational attainment comes from the ACS, and information on average income is calculated from IRS data.

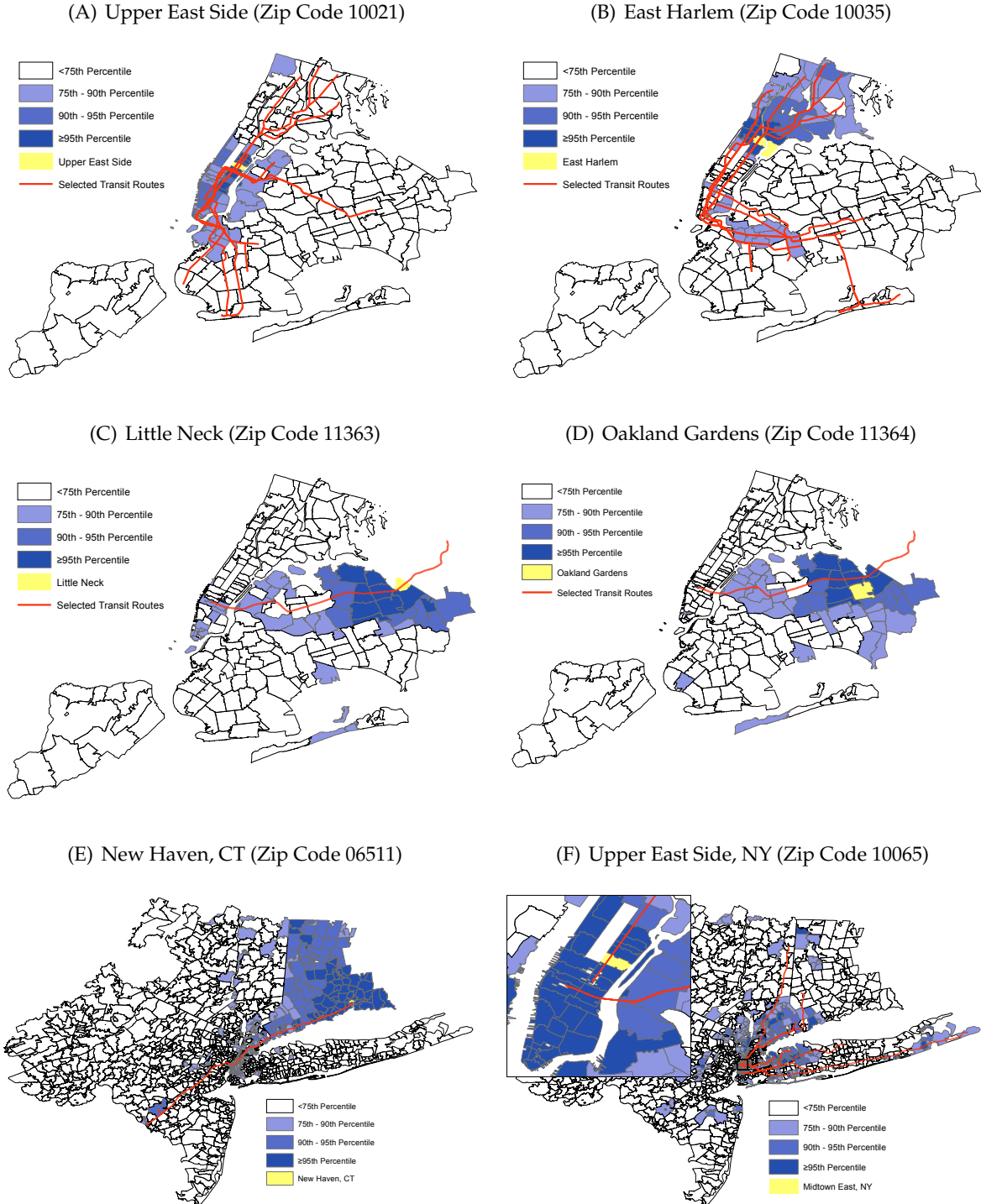
## 2 Determinants of Urban Social Connectedness

**Measuring Social Connectedness.** To compare the intensity of social connectedness between zip codes with varying populations, we construct our measure of *SocialConnectedness<sub>i,j</sub>* as the total number of connections between individuals living in zip code *i* and individuals living in zip code *j*, which we refer to as *FB\_Connections<sub>i,j</sub>*, divided by the product of the number of eligible Facebook users in those zip codes, as in equation 1 (see Bailey et al., 2018b, for the first use of the Social Connectedness Index). This measure represents the relative probability of a Facebook friendship link between a given user in zip code *i* and a given user in zip code *j*:

$$SocialConnectedness_{i,j} = \frac{FB\_Connections_{i,j}}{FB\_Users_i \times FB\_Users_j}. \quad (1)$$

**Social Connectedness in Urban Areas: Case Studies.** Panels A to D of Figure 1 map the percentile ranks of *SocialConnectedness<sub>i,j</sub>* of all zip codes *j* in NYC to four zip codes *i* covering portions of the Upper East Side (10021), East Harlem (10035), Little Neck (11363), and Oakland Gardens (11364), respectively. Relevant transit links are included for illustration. In each panel, relatively more of the connections are mapped to geographically close zip codes; beyond this general pattern, there is substantial heterogeneity in the social networks across the four zip codes. The focal zip codes in panels A and B are roughly two miles apart in uptown Manhattan. The distributions of their respective social networks differ considerably, but essentially all regions with strong social connectedness to these zip codes are linked via direct or one-transfer subway trips. Panel C maps the social network of residents of Little Neck, Queens, a neighborhood on the eastern edge of NYC with easy access to the Long Island Railroad (LIRR) into midtown Manhattan. Little Neck has strong social connectedness to residential areas in midtown Manhattan near the LIRR terminus. Panel D shows the social network of zip code 11364, covering the neighborhood of Oakland Gardens in Queens. While adjacent to Little Neck, which has two LIRR stops, Oakland Gardens does not itself have a LIRR stop. Its

**Figure 1: Social Network Distributions**



**Note:** Figure shows social networks distributions along transit routes. Panels A, B, C, and D show the percentile rank of the relative probability of connection, as measured by  $SocialConnectedness_{i,j}$ , of all zip codes  $j$  in NYC to four zip codes  $i$  in the Upper East Side (Panel A), East Harlem (Panel B), Little Neck (Panel C), and Oakland Gardens (Panel D). Panels E and F show the percentile rank of the relative probability of connection, as measured by  $SocialConnectedness_{i,j}$ , of all zip codes  $j$  in the New York CSA to two zip codes  $i$  in New Haven, CT (Panel E) and the Upper East Side, NY (Panel F). Darker zip codes have a greater probability of connection to a given zip code  $i$ .

social network differs from that of Little Neck in that none of the top connected zip codes extend into Manhattan. The spatial distributions of the social networks presented in Figure 1 therefore provide the first suggestive evidence that NYC’s public transit system plays an important role in enabling the formation and maintenance of social ties across geographic distances. Indeed, it appears as if transit links can effectively “shrink” the geographic distances between locations within the city.

The spatial distribution of social networks of zip codes across the New York CSA also exhibits patterns consistent with those explored for NYC zip codes. Panels E and F of Figure 1 map the percentile rank of  $SocialConnectedness_{i,j}$  for two zip codes  $i$  to all zip codes  $j$  in the New York CSA. Panel E shows the social connectedness to zip code 06511 in New Haven, CT. The social network of New Haven exhibits a strong state border effect along the New York-Connecticut border; it also has a notable instance of long-distance connectivity: over 100 miles away in New Jersey there is a cluster of strongly connected zip codes surrounding the town of Princeton, a feature that is likely driven by students and researchers at Yale University (located in New Haven) and Princeton University; these connections are likely strengthened by the ease of train travel between New Haven and Princeton Junction. Panel F of Figure 1 shows the social network of zip code 10065 in the Upper East Side, which is home to some of the wealthiest residential areas of NYC. This zip code exhibits strong social connectedness to the wealthy northern suburbs as well as to the Hamptons, a popular vacation destination for the well-heeled, roughly 100 miles away at the eastern tip of Long Island. Notably, there are stronger connections to parts of the Hamptons than to zip codes in Long Island City and Astoria, directly across the East River in Queens.

## 2.1 Geographic Distance, Social Distance, and Social Connectedness

The previous section presented a number of case studies that suggest a relationship between social connectedness, geographic distance, transit availability, and demographic similarity. We next estimate the elasticity of social connectedness with respect to these objects more formally.

To systematically measure the ease of travel between two zip codes, we use the Google Maps API to collect travel times on public transit between the geographic centers of all zip codes on a weekday morning,<sup>3</sup> and measure cab cost in dollars using data from the New York Taxi and Limousine Commission (TLC).<sup>4</sup> There is substantial variation in transit travel time and cab costs over similar

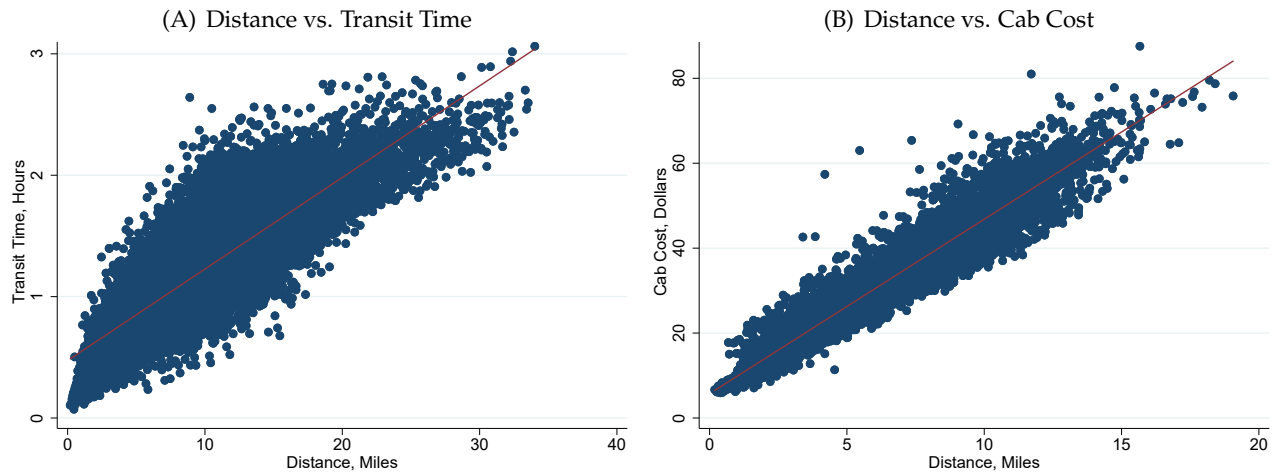
<sup>3</sup>Our data on public transit time is collected from the Google Distance Matrix API. For each pair of zip codes  $i$  and  $j$  we collected the transit time of a trip from  $i$  to  $j$  and from  $j$  to  $i$ . Trips originate and end at the geographic centers of each zip code. The transit time measure between two zip codes ( $TravelTime_{i,j}$ ) is the average time of trip  $i$  to  $j$  and trip  $j$  to  $i$ . We queried Google for travel times on a weekday, March 15th, 2017. We queried travel times more than two weeks in advance, so that contemporaneous delays or construction work not would influence trip times. We pulled the travel time on a weekday morning for a traveler that has to arrive at the other zip code by 9AM, to estimate travel time on a work day.

<sup>4</sup>The TLC reports the data for each cab trip taken in the first six months of 2016. The latitude and longitude of the origin and destination of 19.7 million trips, composed of 11.2 million yellow cab trips and 8.5 million green cab (“borough cab”) trips, were matched to their origin and destination zip codes. For green cabs, which primarily serve the outer boroughs, all trips taken in the first six months of 2016 were matched to zip codes. For yellow cabs, which provide a greater share of trips but are more concentrated in Manhattan, only trips in March, 2016, were matched to zip codes. The cost of a trip from zip code  $i$  to zip code  $j$  is calculated as the average of the costs of all trips originating in zip code  $i$  and ending in zip code  $j$ . We only consider zip codes that have at least one trip in each direction, and calculate the cost of travel between a zip code-pair composed of zip codes  $i$  and  $j$  as the average of the cost of trips from  $i$  to  $j$  and trips from  $j$  to  $i$ .

geographic distances in NYC. For example, Figure 2 shows that the 95th percentile transit trip time between zip codes that are (roughly) 2.5 miles apart is only 4 minutes less than the 5th percentile trip time between zip codes that are 10 miles apart. Much of this variation is driven by a combination of public transit infrastructure and geography. For example, a transit trip between the East Village and Greenpoint, two neighborhoods facing one another across the East River that lack a connection via a tunnel or bridge, is at the 90th percentile of trip time compared to trips between other zip code-pairs separated by a similar geographic distance. Figure 2 also shows there is substantial variation of cab costs for similar distances. As an alternative way of presenting the same information, Panel A shows a scatter plot demonstrating the variation of transit travel times between zip code pairs that are a given distance apart from each other; Panel B shows the variation in cab costs for zip code pairs that are a given distance apart.

**Figure 2: Travel Time and Costs Variation by Distance, NYC**

	Transit Travel Times Between Zips That Are:			Cab Trip Costs Between Zips That Are:		
	2.5 Miles Apart	5 Miles Apart	10 Miles Apart	2.5 Miles Apart	5 Miles Apart	10 Miles Apart
Mean	0:33:51	0:51:16	1:18:45	\$15.61	\$26.23	\$47.40
P5	0:18:45	0:32:46	0:53:44	\$12.39	\$20.83	\$39.12
P10	0:20:51	0:34:06	0:56:45	\$13.10	\$22.08	\$41.52
P25	0:25:53	0:41:23	1:07:49	\$13.71	\$23.96	\$44.48
Median	0:33:44	0:49:51	1:17:51	\$14.82	\$25.81	\$47.23
P75	0:39:43	0:59:12	1:28:22	\$16.42	\$28.44	\$50.63
P90	0:47:02	1:08:43	1:40:17	\$18.59	\$31.03	\$53.56
P95	0:49:59	1:17:49	1:47:04	\$23.50	\$32.16	\$54.73
N	138	224	225	131	194	94



**Note:** Table shows across-zip-code-pair summary statistics for transit time and cab trip cost between zip codes in a zip code-pair at various distances. All travel times and cab costs for zip code-pairs that are between  $\pm 1$  miles of the indicated distance are included in each column. Not all zip code pairs were traveled via cab during our sample period, allowing us to only calculated cab trip costs for a subset of zip code pairs. We also show scatter plots at the zip code-pair level of transit time (Panel A) and cab trip cost (Panel B) on the vertical axes. The horizontal axes for both panels show the geographic distance between the centers of each zip code pair.

To obtain a more systematic understanding of the effect of transportation links on social networks, we next use equation 2 to explore the pairwise friendship links between zip codes:

$$\log(\text{SocialConnectness}_{ij}) = \beta_0 + \beta_1 \log(d_{ij}) + X_{ij} + \psi_i + \zeta_j + \epsilon_{ij}. \quad (2)$$

The dependent variable is the log of social connectedness (defined in equation 1), and  $\log(d_{ij})$  denotes the log of the “distance” between  $i$  and  $j$ . Here, “distance” will be variously defined as the geographic distance between the central points of zip codes  $i$  and  $j$ , the public transit time between the central points of zip codes  $i$  and  $j$ , and the average cost of cab trips between zip codes  $i$  and  $j$ . Control variables  $X_{ij}$  include measures of the dissimilarity of the two zip codes along demographic and socioeconomic factors. These factors are income (the difference in average income across the zip code-pair), education (the difference in the shares of residents without a high school degree across the zip code-pair), and race (the difference in the non-Hispanic white shares of the populations across the zip code-pair). All specifications include fixed effects  $\psi_i$  and  $\zeta_j$  for zip codes  $i$  and  $j$ , respectively.

Table 1 shows the results of the regression 2 with  $\log(d_{ij})$  representing geographic distance in columns 1 and 2, public transit time in columns 3 and 4, and cab cost in columns 5 and 6. Columns 1, 3, and 5 are the baseline specifications as shown in regression 2. Columns 2, 4, and 6 include interaction terms for pairs of zip codes that are both in the top third of the income distribution and pairs that are both in the bottom third of the income distribution with  $\log(d_{ij})$ , to test if the social connectedness of zip codes responds differently to transit times or cab costs based on differences in zip code incomes. When we compare columns 1 and 3, we find that the coefficient for transit time is over 60% greater in magnitude than that for geographic distance. The estimates imply that a 10% greater geographic distance between zip codes is associated with 8.7% lower social connectedness, while a 10% increase in public transit time is associated with 14.2% lower social connectedness. Likewise, column 5 indicates that a 10% increase in cab cost is associated with a 10.6% decline in social connectedness. These results suggest that public transportation infrastructure plays a more important role in the formation of social networks in urban settings than simple geographic distance does.

Table 1 also documents that, beyond the various measures of distance, zip codes that are more similar in terms of their education levels and their racial composition are more likely to be socially connected, providing evidence for homophily within New York City. For example, conditional on the geographic distance and differences in income and education levels, a 10 percentage point increase in the difference in the share of the population that is white is associated with about a 11% to 12% decline in social connectedness. Similarly, a 10 percentage point increase in the difference of the population shares with no high school is associated with a 7%-10% decline in social connectedness. While differences in income do not imply differences in social connectedness (once we condition for differences in racial composition and educational attainment), we do find that the elasticity of social connectedness to the various measures of distance is larger when zip codes have very different income measures. In particular, columns 2, 4, and 6 show that the effect of increasing distance on



**Table 1:** The Effect of Distance and Transportation on Social Connectedness, NYC

	(1)	(2)	(3)	(4)	(5)	(6)
Log(Distance in Miles)	-0.872*** (0.044)	-0.951*** (0.044)				
Log(Avg. Time on Transit)			-1.418*** (0.067)	-1.498*** (0.071)		
Log(Avg. Cab Cost)					-1.059*** (0.044)	-1.113*** (0.048)
$\Delta$ Share Pop White (%)	-0.012*** (0.001)	-0.011*** (0.001)	-0.013*** (0.001)	-0.012*** (0.001)	-0.011*** (0.001)	-0.011*** (0.001)
$\Delta$ Share Pop No High School (%)	-0.010*** (0.002)	-0.007*** (0.002)	-0.009*** (0.002)	-0.007*** (0.002)	-0.008*** (0.003)	-0.006** (0.002)
$\Delta$ Avg. Income (k\$)	-0.000 (0.000)	-0.001 (0.000)	-0.001 (0.001)	-0.000 (0.000)	-0.001 (0.000)	-0.001** (0.000)
Interaction: Rich Zip-Pair (X Dist., Transit Time, or Cab Cost)		0.234*** (0.064)		0.150* (0.083)		0.444*** (0.072)
Interaction: Poor Zip-Pair (X Dist., Transit Time, or Cab Cost)		0.199*** (0.050)		0.274*** (0.080)		0.020 (0.064)
Dummy for Zip-Pair Type		Y		Y		Y
Zip Code Fixed Effects	Y	Y	Y	Y	Y	Y
Number of Observations	16,283	16,283	16,283	16,283	7,873	7,873
R-Squared	0.759	0.765	0.759	0.763	0.836	0.841

**Note:** Table shows results from regression 2. The unit of observation is a zip code-pair. The dependent variable in all columns is the log of  $SocialConnectedness_{i,j}$  as defined in equation 1. All specifications include zip code fixed effects and measures of the similarity of zip codes within the pair along socioeconomic and demographic dimensions. The measure of “distance” in regression 2 is variously defined as geographic distance (columns 1-2), transit time (columns 3-4), and cab cost (columns 5-6). Columns 2, 4, and 6 include interaction terms for rich zip code-pairs and poor zip code-pairs with “distance.” Coefficients for the dummy variables are excluded for brevity. Standard errors are double clustered by each zip code  $i$  and zip code  $j$  in a zip code-pair. Significance levels: \* ( $p < 0.10$ ), \*\* ( $p < 0.05$ ), \*\*\* ( $p < 0.01$ ).

social connectedness is smaller across zip code pairs with similar incomes (i.e., zip code pairs where both zip codes are in the top tertile or those where zip codes are in the bottom tertile of the income distribution). Said differently, reducing travel times appears to have a disproportionate effect on fostering social connectedness across regions with very different incomes.

In order to examine how the distance affects social connectedness at the CSA level, Table 2 shows the result from performing regression 2 for zip codes across the New York CSA. Since many of these zip codes are not well connected via public transport, we use the log of geographic distance as the measure of  $\log(d_{ij})$ . Column 1 excludes zip code fixed effects and socioeconomic dissimilarity variables  $X_{i,j}$ , and column 2 includes zip code fixed effects but excludes socioeconomic dissimilarity variables  $X_{i,j}$ . Column 3 includes an additional variable indicating whether both zip codes are within the same state. Column 4 includes differences in demographic variables, and column 5 adds interac-

**Table 2:** The Effect of Geographic Distance on Social Connectedness, New York CSA

	(1)	(2)	(3)	(4)	(5)
Log(Distance in Miles)	-1.229*** (0.026)	-1.582*** (0.023)	-1.383*** (0.024)	-1.268*** (0.025)	-1.329*** (0.024)
Same State			1.081*** (0.040)	1.158*** (0.039)	1.155*** (0.039)
$\Delta$ Share Pop White (%)				-0.010*** (0.001)	-0.010*** (0.001)
$\Delta$ Share Pop No High School (%)				-0.021*** (0.002)	-0.017*** (0.002)
$\Delta$ Avg. Income (k\$)				-0.006*** (0.000)	-0.005*** (0.000)
Interaction: Rich Zip-Pair (X Dist. Miles or Transit Time)					0.183*** (0.030)
Interaction: Poor Zip-Pair (X Dist. Miles or Transit Time)					0.196*** (0.032)
Dummy for Zip-Pair Type					Y
Zip Code Fixed Effects		Y	Y	Y	Y
Number of Observations	625,743	625,741	625,741	625,741	625,741
R-Squared	0.389	0.714	0.738	0.784	0.788

**Note:** Table shows results from regression 2 for zip code-pairs in the New York CSA. The unit of observation is a zip code pair. The dependent variable in all columns is the log of  $SocialConnectedness_{i,j}$  as defined in equation 1. The measure of “distance” is geographic distance in all specifications. Column 1 does not include zip code fixed effects and controls. Column 2 includes zip code fixed effects. Column 3 incorporates a control variable for zip codes that are in the same state. Column 4 adds measures of the similarity of zip codes along socioeconomic and demographic dimensions. Column 5 additionally includes interaction terms for rich zip codes and poor zip codes. Coefficients for the dummy variables for the various zip pair types are excluded for brevity. Standard errors are double clustered by each zip code  $i$  and zip code  $j$  in a zip code-pair. Significance levels: \* ( $p < 0.10$ ), \*\* ( $p < 0.05$ ), \*\*\* ( $p < 0.01$ ).

tion terms for rich zip code-pairs and poor zip code-pairs, defined as above. The effect of distance in all specifications is greater for zip codes across the CSA than it is for the subset of zip codes within NYC. This is consistent with prior research demonstrating that urban social networks are less geographically determined than those over larger areas (Herrera-Yague et al., 2015). The coefficients on distance in these regressions are generally smaller in magnitude than the ones for regressions in earlier research by Bailey et al. (2018a) at the county level for counties within 200 miles of one another (this is the relevant comparison, as there are very few zip code-pairs more than 200 miles apart in the New York CSA). This difference may be due to differences in the properties of social networks measured at this level of aggregation, or due to our sample of zip codes centered on a large urban area where the effect of distance is weaker.

Column 3 of Table 2 shows that the social connectedness between two zip codes in the same state is about twice as large as the connectedness between equidistant zip codes across different states.



This could, for example, be the result of school districts that do not cross state lines; other possible explanations include the role of occupational licensing in restricting cross-state moves, and thereby cross-state friendship formation. The negative coefficients on all of the socioeconomic dissimilarity measures in column 4 are suggestive of homophily in the New York CSA; homophily based on income and educational attainment seems stronger in the New York CSA relative to NYC, while homophily based on race appears similarly large. Finally, the estimates in column 5 support the findings from the within-NYC analysis: the social connectedness across zip codes with populations of different income drops off faster in distance than the social connectedness across zip codes with more similar incomes.

### 3 The Geographic Concentration of Social Networks

In this section, we document heterogeneity in the geographic concentration of social networks across zip codes. We also explore which factors are associated with the geographic dispersion of these networks. In Section 3.1, we explore heterogeneity in two different measures of social network concentration. In Section 3.2, we investigate the relationship between the geographic dispersion of social networks and socioeconomic outcomes such as income and education. In Section 3.3, we analyze the relationship between the concentration of social networks and the ease of access to public transit.

#### 3.1 Measurement of Social Network Concentration

We consider two measures of the geographic concentration of social networks: the share of friends that lives within a certain geographic radius (e.g., 1 mile or 5 miles), and the share of friends that lives within a certain number of people (e.g., within the nearest 1 million or 5 million people).

To construct our concentration measures for small distances such as one mile, we have to determine which friends are included within this range, even though we only observe the locations of individuals and their friends at the zip code level. We therefore construct our measures by weighting friendships to individuals in each region  $j$  by the population-weighted share of census blocks in region  $j$  that are within that distance of the population-weighted center of zip code  $i$ . Specifically, we use the following equation to construct our measure of the geographic concentration of zip code  $i$ 's friendship network:

$$ShareWithinDMiles_i = \sum_j ShareFriends_{i,j} * \frac{\sum_{j_b} Pop_{j_b} * \mathbb{1}_{d_{i,j_b} \leq D}}{TotalPop_j} \quad (3)$$

Here,  $d_{i,j_b}$  indicates the distance from the population-weighted center of zip code  $i$  to the center of each census block  $j_b$  in region  $j$ . We find the population of each region  $j$  that is within a given distance  $D$  from zip code  $i$  by summing the population of all census blocks  $j_b$  for which  $d_{i,j_b}$  is less than  $D$ , and divide this by the total population of region  $j$ . We then weigh the share of friends of zip code  $i$  living in region  $j$ , given by  $ShareFriends_{i,j}$ , by the share of the population of zip code  $j$  that lives within  $D$  miles of the center of zip code  $i$ , before summing over all regions  $j$ . We will use the following two

objects as our measures of the geographic concentration of social networks. For our first measure, the share of friends living within a certain radius,  $D$  represents one, five, ten, or fifty miles. For our second measure, the share of friends living within a certain number of people, we define  $D$  as the radius from the center of each zip code  $i$  that contains a given number of people, and then construct the statistics as above based on that distance.

**Table 3:** Summary Statistics of Geographic Concentration of Social Networks

(A) NYC

	Share of Friends Living Within:			Share of Friends Among Nearest:		
	1 Mile	5 Miles	10 Miles	250K People	1 Mil. People	10 Mil. People
Mean	6.3%	29.3%	44.0%	10.2%	21.9%	55.1%
P5	3.3%	19.5%	34.1%	4.9%	13.1%	38.8%
P10	4.0%	22.2%	36.0%	5.9%	14.2%	40.2%
P25	4.8%	26.2%	39.5%	7.3%	17.4%	49.3%
Median	6.2%	29.0%	44.0%	9.5%	22.0%	58.2%
P75	7.4%	32.4%	48.4%	11.6%	25.4%	61.2%
P90	8.8%	37.2%	52.7%	14.4%	29.5%	64.2%
P95	10.0%	39.6%	53.8%	19.6%	32.7%	66.0%

(B) New York CSA

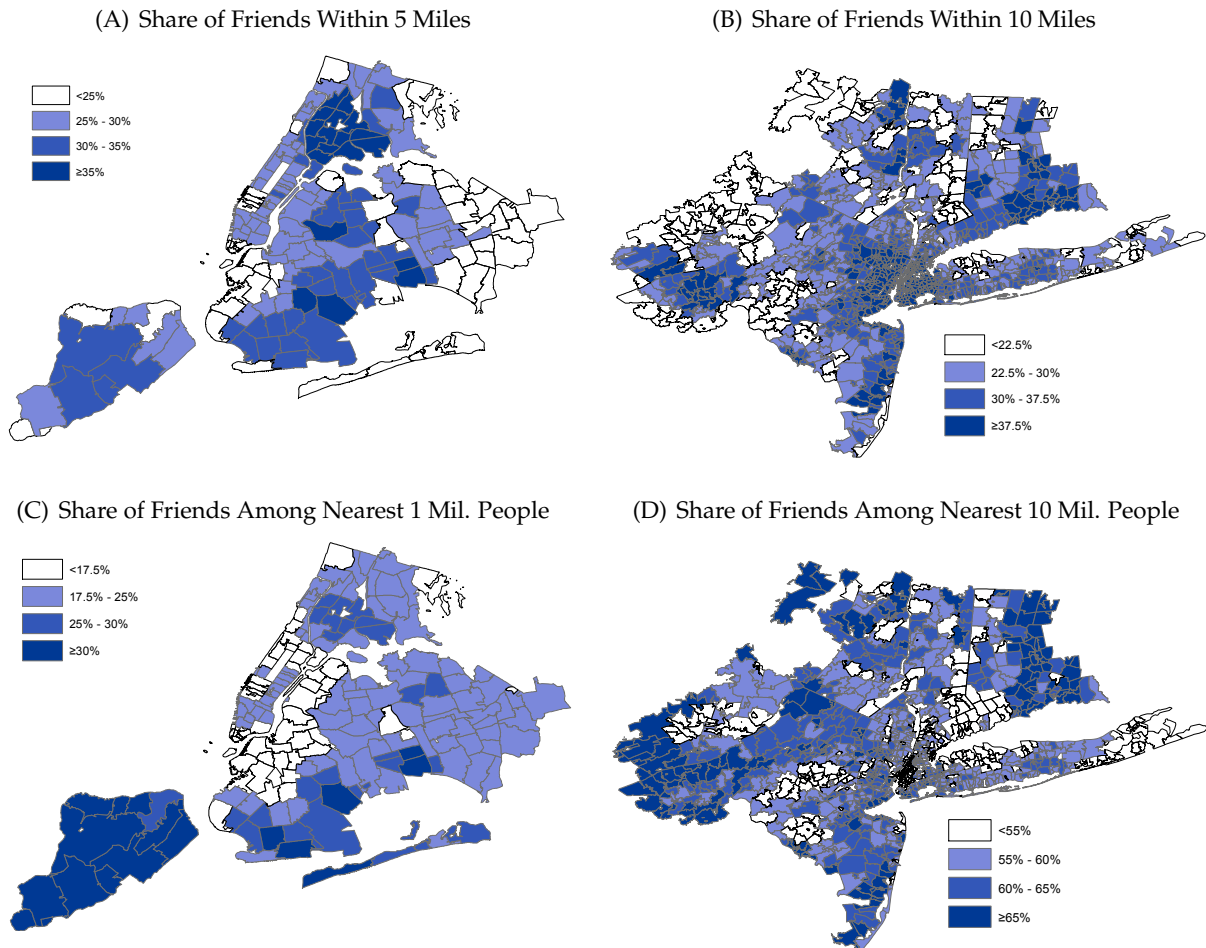
	Share of Friends Living Within:			Share of Friends Among Nearest:		
	5 Miles	10 Miles	50 Miles	1 Mil. People	10 Mil. People	50 Mil. People
Mean	25.8%	38.1%	64.1%	32.3%	57.9%	75.4%
P5	12.2%	22.4%	48.6%	14.8%	40.6%	61.5%
P10	14.8%	25.1%	53.0%	17.5%	47.3%	65.8%
P25	20.2%	31.1%	60.0%	22.9%	54.6%	73.8%
Median	25.8%	38.4%	65.9%	31.4%	59.1%	76.8%
P75	31.0%	45.7%	69.1%	41.4%	62.4%	79.2%
P90	37.1%	51.1%	72.0%	49.2%	66.1%	81.1%
P95	40.4%	53.7%	73.4%	52.3%	68.5%	81.9%

**Note:** Table shows summary statistics of the geographic concentration of social networks. Panel A shows across-zip code summary statistics of the share of domestic friends of a zip code's population that live within 1, 5, and 10 miles of a zip code, and the share of domestic friends of a zip code's population that are among the nearest 250 thousand, 1 million, and 10 million people in and surrounding a zip code for NYC. Panel B shows across-zip code summary statistics of the share of domestic friends of a zip code's population that live within 5, 10, and 50 miles of a zip code, and the share of domestic friends of a zip code's population that are among the nearest 1 million, 10 million, and 50 million people in and surrounding a zip code for the New York CSA. Zip codes are weighted by their populations.

Panel A of Table 3 provides summary statistics at the zip code level of the geographic concentration of social networks in NYC, based on the distribution of U.S. Facebook friends of users residing in each zip code. For the residents of the population-weighted average zip code in NYC, 6.3% of U.S. friends live within one mile, 29.3% of U.S. friends live within five miles, and 44.0% of U.S. friends live within ten miles. There is significant heterogeneity in the geographic concentration of friendship

links: across zip codes, the 5-95 percentile range of U.S. friends living within one mile is 3.3% to 10.0%. Similarly, Panel B of Table 3 provides summary statistics on the geographic concentration of U.S. friendship networks by zip code for the New York CSA. For the residents of the population-weighted average zip code in the New York CSA, 25.8% of U.S. friends live within five miles, 38.1% of U.S. friends live within ten miles, and 64.1% of U.S. friends live within fifty miles. Once again, there is significant heterogeneity in the concentration of friendship connections: the 5-95 percentile range of friends living within ten miles is 22.4% to 53.7%.

**Figure 3: Geographic Concentration of Social Networks**



**Note:** Figure shows the geographic concentration of social networks for each zip code in New York. Panel A shows a map at the zip code level of the share of all U.S. friends that live within 5 miles for each NYC zip code. Panel B shows a map of the share of all U.S. friends that live within 10 miles for each zip code in the New York CSA. Panel C shows a map at the zip code level of the share of all U.S. friends that are among the nearest 1 million people for each NYC zip code. Panel D shows a map of the share of all U.S. friends that are among the nearest 10 million people for each zip code in the New York CSA.

Panel A of Figure 3 maps the spatial distribution of the share of U.S. friends living within five miles for each zip code in NYC. Zip codes with the most geographically dispersed friendship networks are primarily in the western area of Brooklyn and the eastern portion of Queens, as well as the Downtown

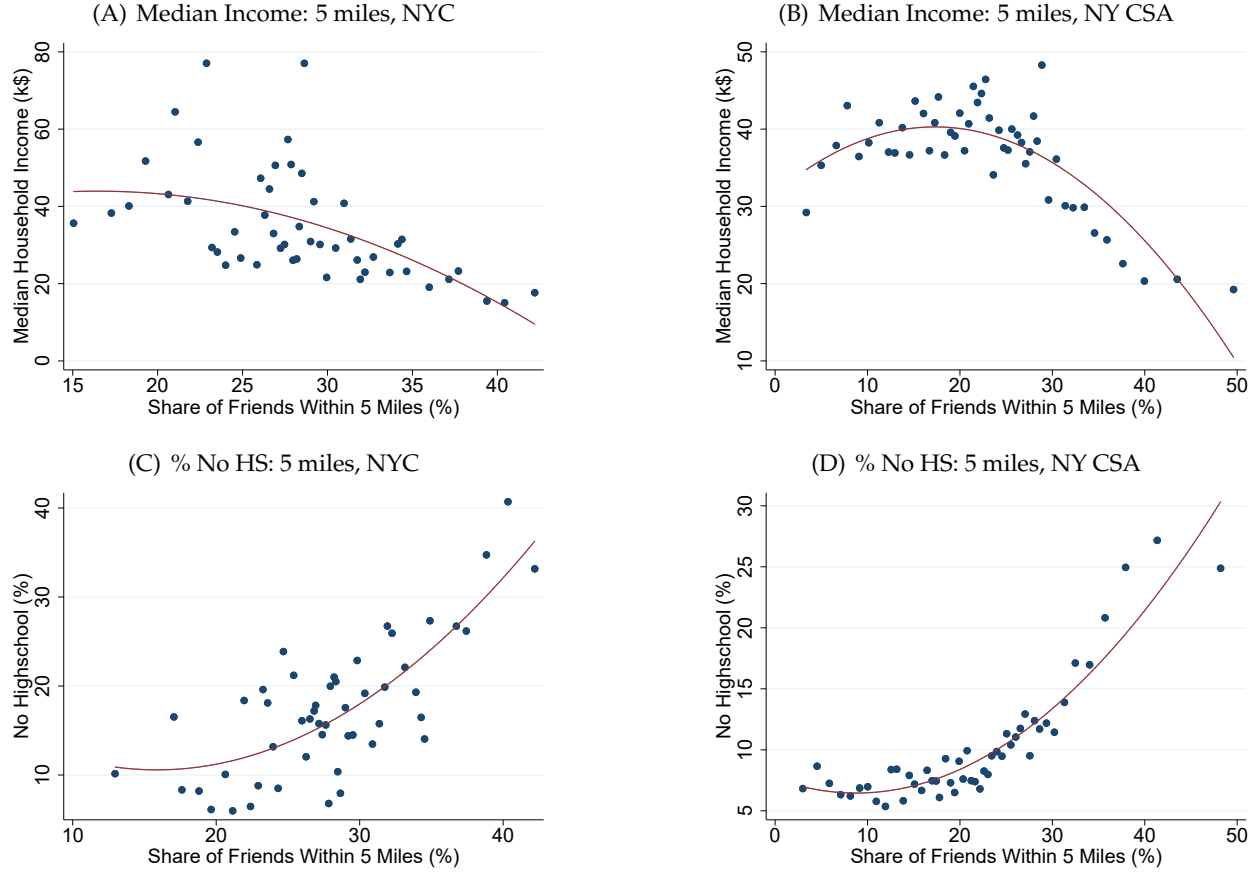
and Midtown West neighborhoods of Manhattan. Panel B of Figure 3 shows the spatial distribution of the share of friends within ten miles for zip codes within the New York CSA, revealing that networks are generally more geographically concentrated in the urban areas within the CSA, with high concentrations most evident in the area in and surrounding NYC but also present in New Haven, CT, Allentown, PA, and Seaside Heights, NJ.

We find similar heterogeneity in the share of friends living within a certain number of people: Panel A of Table 3 indicates that for population-weighted average zip code, 21.9% of friendship links are to the one million closest individuals, but this number ranges from 13.1% to 32.7% between the 5th and the 95th percentiles of the zip code distribution. Panel B of Table 3 also highlights that for the residents of the population-weighted average zip code in the New York CSA, 32.3% of U.S. friends are among the nearest one million people, 57.9% of U.S. friends are among the nearest ten million people, and 75.4% of U.S. friends are among the nearest fifty million people. There is also a high degree of heterogeneity for these measures: the 5-95 percentile range of friends living among the nearest one million people is 14.8% to 52.3%. Panel C of Figure 3 maps the share of friends of individuals who live within the nearest one million people for NYC zip codes. Notably, the aggregate social networks of users in population-dense regions, such as those in north Brooklyn, are comparatively more dispersed in terms of the share of friends within a certain number of people than in terms of the share of friends living within a certain geographic distance. The opposite pattern characterizes the social networks of Staten Island, the NYC borough with the lowest population density. Panel D of Figure 3 shows the spatial distribution of social network density in the New York CSA, using as the measure the share of friends among the nearest 10 million people. The distribution is different from that in Panel B, as the urban cores and inner suburbs display less social network density using this measure. The differences are primarily driven by variation in population densities across urban and non-urban areas within the New York CSA. It also suggests that physical distance to core urban areas may not be the only determinant of social connectedness, and echoes with findings from Section 2 that other factors such as transit infrastructure may also matter.

### 3.2 Socioeconomic Outcomes and the Concentration of Social Networks

We next explore the geographic concentration of social networks is correlated with observable individual characteristics at the zip code level. Figure 4 shows zip code-level binned scatter plots of the share of friends living within 5 miles against income and education measures in NYC (left panel) and the New York CSA (right panel); similar patterns arise when we measure the concentration of social networks at other distances or as the share of friends within a certain number of people. The binned scatter plots illustrate that zip codes with more widely dispersed social networks generally have higher incomes and education levels. While the relationships in Figure 4 are not necessarily causal, the literature has proposed many causal mechanisms for the observed patterns: indeed, access to diverse information through broad social networks is central to many theories of innovation, social mobility, and economic growth (Jackson, 2014; Granovetter, 2005).

**Figure 4: Demographics and Geographic Concentration of Social Networks**



**Note:** Figure shows binned scatter plots with zip codes as the unit of observation. The left column includes all NYC zip codes, while the right column includes all zip codes in the NY CSA. The horizontal axis plots the share of the U.S.-based friends of a zip code's population that live within 5 miles of a zip code. Row 1 plots the median income for residents of a zip code on the vertical axis, and row 2 the share of the zip code's population without a high school degree. The R-Squared values corresponding to the quadratic line of fit are: 15.0% (Panel A), 15.8% (Panel B), 27.9% (Panel C), 32.8% (Panel D).

### 3.3 Ease of Transit and the Concentration of Social Networks

Having established that there is substantial heterogeneity in the geographic concentration of social networks, we next explore whether differences in the public transit infrastructure across zip codes can explain this heterogeneity. We construct two measures of the ease of public transit at the zip code level, which we call “transit inconvenience” and “transit access.” Transit access is measured as the share of the zip code's population that lives within a quarter mile of a rail transit station.<sup>5</sup> Transit inconvenience is based on the travel times computed in Section 2, and constructed as the average of

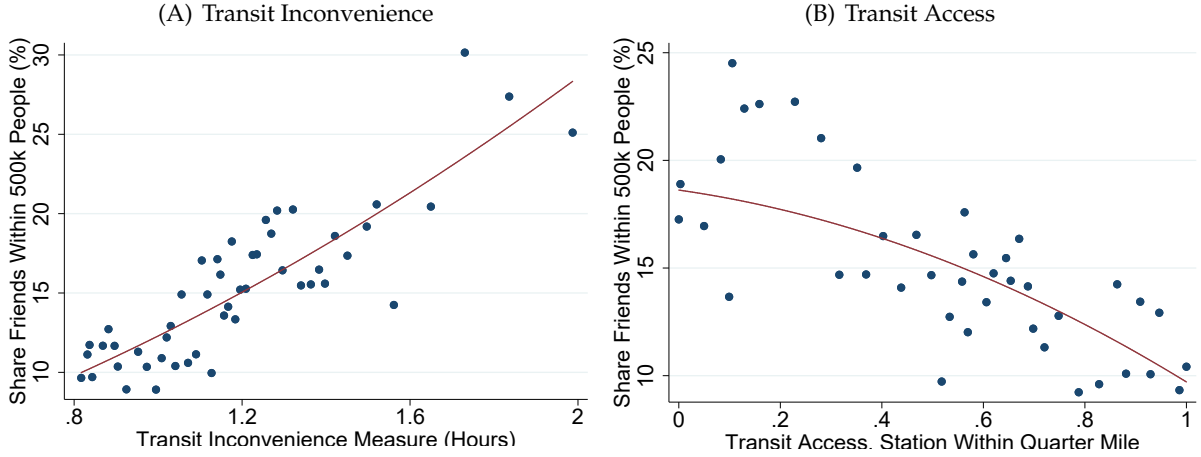
<sup>5</sup>A rail transit station is defined as either an MTA subway stop or a Long Island Railroad (LIRR) stop, as these are the two most important rail transit options within NYC. This transit access measure is intended to capture access to physical rapid transit infrastructure. Of course, zip codes may have access to other forms of public transit, and the measure of public transit time that we collect from Google allows for transit via any vehicle (trains as well as buses, ferries, trams, etc.), but rail transit provides the majority of public transit trips within the city (MTA, 2016a,b).

$TravelTime_{i,j}$  for each zip code  $i$  with all zip codes  $j$  over the number of zip code observations  $n_j$ :<sup>6</sup>

$$TransitInconvenience_i = \frac{\sum_j TravelTime_{i,j}}{n_j}. \quad (4)$$

Panel A of Figure 5 shows a binned scatter plot of the relationship between the inconvenience of transit measured in hours, as defined in equation 4, and the geographic concentration of social networks measured by the share of friends who live within the nearest 500k people for zip codes in NYC. The social networks of those zip codes with more convenient transit are less geographically concentrated compared to those with less convenient transit. Panel B of Figure 5 shows a binned scatter plot relating the transit access of zip codes to the share of friends living within the nearest 500k people. In this case, zip codes with greater access to public transit have more geographically dispersed social networks. Overall, the findings from these plots are consistent with the notion that the ease of transportation is associated with a wider geographic dispersion of social networks.

**Figure 5:** Ease of Transit and Social Network Concentration, NYC



**Note:** Figure shows binned scatter plots at the zip code level of the share of friends within the nearest 500k people on the vertical axes. The horizontal axis in Panel A shows the average transit time from a zip code to all other zip codes within New York City as defined in equation 4. The horizontal axis in Panel B shows the share of the population of each zip code that is within a quarter mile of a subway or LIRR stop.

There are many potentially confounding factors that could influence the relationship between the share of friends within certain distance or number of people and the ease of travel via public transit. For instance, due to the radial design of New York’s subway system, all but one train service run through the relatively wealthy areas of midtown or downtown Manhattan. To separately explore the role of transit infrastructure beyond the demographic measures that it is correlated with, we next estimate regression 5:

$$ShareWithin500k_i = \beta_0 + \beta_1 Transit_i + \beta_2 X_i + \epsilon_i \quad (5)$$

<sup>6</sup>All results in this section are similar for measures of transit inconvenience that weight each zip code  $j$  in equation 4 by the population in zip code  $j$ , so that having a longer travel time to a high-population zip code counts more towards transit inconvenience than high travel time to a low-population zip code does.

The dependent variable is the geographic concentration of social networks, measured as the share of friends that live within the nearest 500k people, though our conclusions are similar when using our other measures of the geographic concentration of social networks. Depending on the specification,  $Transit_i$  will represent the transit inconvenience measure (equation 4), or the share of a zip codes' population within a quarter mile of a transit stop.  $X_i$  includes controls for socioeconomic characteristics of each zip code.

**Table 4:** Transit and the Geographic Dispersion of Social Networks, NYC

	(1)	(2)	(3)	(4)	(5)	(6)
Transit Inconvenience (Hours)		14.844*** (1.260)		15.321*** (1.476)		13.881*** (1.899)
Share Pop 1/4 Mile from Transit (%)			-0.085*** (0.012)		-0.083*** (0.013)	-0.017*** (0.014)
Share Pop White (%)	0.043* (0.024)			0.010 (0.019)	0.051** (0.021)	0.015 (0.019)
Share Pop No High School (%)	0.005 (0.062)			0.087* (0.050)	0.131** (0.059)	0.106** (0.052)
Avg. Income (k\$)	-0.017*** (0.004)			0.002 (0.004)	-0.008** (0.004)	0.003 (0.004)
Number of Observations	182	182	182	182	182	182
R-Squared	0.110	0.435	0.224	0.447	0.285	0.451

**Note:** Table shows the results from regression 5. The unit of observation is a NYC zip code. The dependent variable is the share of friends that live within the nearest 500k people. Significance levels: \* ( $p < 0.10$ ), \*\* ( $p < 0.05$ ), \*\*\* ( $p < 0.01$ ).

The estimates from regression 5 are presented in Table 4. Column 1 shows that differences in demographic and socioeconomic characteristics explain only about 11% of the across-zip code variation in the geographic concentration of social networks. Columns 2 and 3 confirm that there is a statistically significant positive relationship between the transit inconvenience measure and the geographic concentration of social networks and a statistically significant negative relationship between transit access and the geographic concentration of the networks, respectively. The R-Squared measures in these univariate regressions are substantially larger than those for the regression presented in column 1, suggesting that much of the differences in social network concentration as measured here are explained by differences in the ease of public transit. Quantitatively, a 15 minute increase in the average travel time to all zip codes is associated with a 3.7 percentage point increase in the share of friends living within the nearest 500k people. Similarly, a ten percentage point increase in the zip code population that lives within a quarter mile of a transit stop is associated with a 0.8 percentage point decline in the share of friends living within the nearest 500k people. Columns 4 and 5 show that the estimated relationship between social network concentration and the ease of public transit is unaffected by the addition of the demographic and socioeconomic controls. Finally, in column 5 we include controls for both of our measures of ease of transit. The primary variable that explains



variation in social network concentration in this specification is the average public transit travel time required to travel to NYC zip codes.

Overall, the results in this section can be summarized as follows. First, there is substantial heterogeneity across zip codes in various measures of the geographic concentrations of their social networks. Second, zip codes with more concentrated social networks generally perform worse on socioeconomic indicators such as income and education levels. Third, much of this variation in social network concentration is explained, at least statistically, by variation in the ease of travel via public transit to the rest of NYC. These results are highly consistent with stories in which investments in public transportation infrastructure allow individuals to form and maintain more geographically dispersed networks, which can expose those individuals to a more diverse set of ideas and opportunities, and thereby contribute to agglomeration externalities.

## 4 Connected Communities in New York

We next provide an alternative description of the geographic structure of social networks across New York. To do this, we use a hierarchical agglomerative linkage clustering algorithm to construct hypothetical “communities” of zip codes that maximize within-group social connectedness. This procedure allows us to determine which groups of zip codes groups are maximally connected to one another, and to compare the resulting connected communities to existing administrative boundaries, such as NYC boroughs or states.

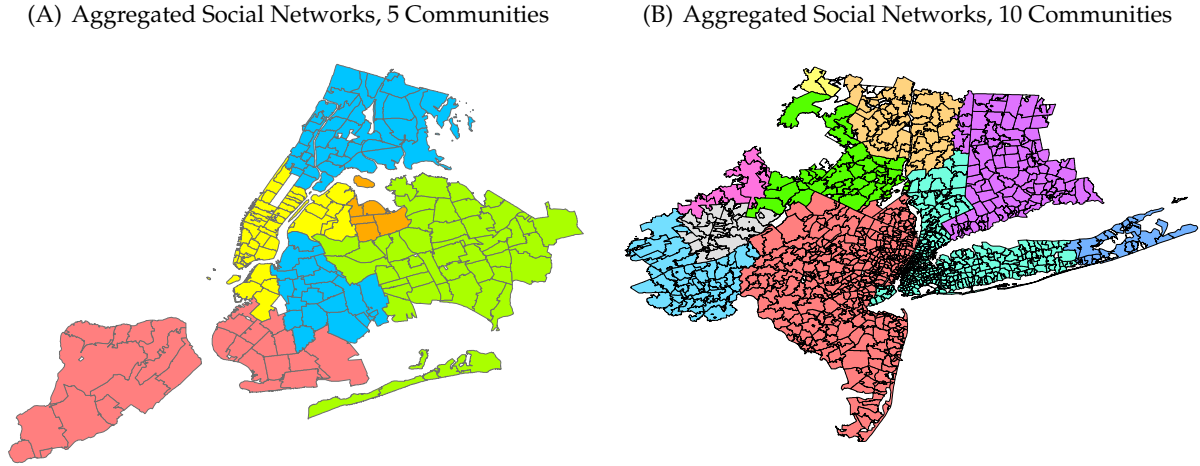
The algorithm starts by considering each of the  $N$  zip codes in a region (either NYC or the New York CSA) as separate communities of size one. The two “closest” zip codes, based on their relationships with all other zip codes, are then merged into one larger community, thus producing  $N - 1$  total communities. We define the “distance” between two zip codes as the inverse of  $SocialConnectedness_{i,j}$  in equation 1. The “distance” between the newly formed community  $i$  and each other zip code  $j$  is then calculated as the average of the “distances” for both of the constituent zip codes in the community to each zip code  $j$ . The two most connected communities are then again merged, producing  $N - 2$  total communities. This process continues until all zip codes are merged into a given number of “connected communities.”

Panel A of Figure 6 shows the result of grouping NYC zip codes into five connected communities. A large band of Brooklyn is clustered together with Harlem and the Bronx; interestingly, this connected community thus consists of two non-contiguous elements that are more connected with each other than they are with Manhattan, which lies between them. This finding again suggests that geographic distance might not be as relevant a measure of “distance” within dense urban areas as it is at other levels of aggregation. Manhattan below Harlem and Morningside Heights joins with a handful of neighborhoods across the East River in Brooklyn and Queens; Brooklyn south of Prospect Park to Coney Island is grouped with Staten Island; and the rest of Queens is split into a small northern community adjacent to LaGuardia Airport and a large eastern community.

We also repeat the hierarchical agglomerative clustering for all zip codes in the New York CSA.



**Figure 6:** Agglomerative Linkage Clustering of Communities



**Note:** Figure shows the results of the hierarchical agglomerative linkage clustering algorithm. Panel A shows NYC zip codes grouped together to create 5 connected “communities.” Panel B shows New York CSA zip codes grouped together to create 10 “communities.”

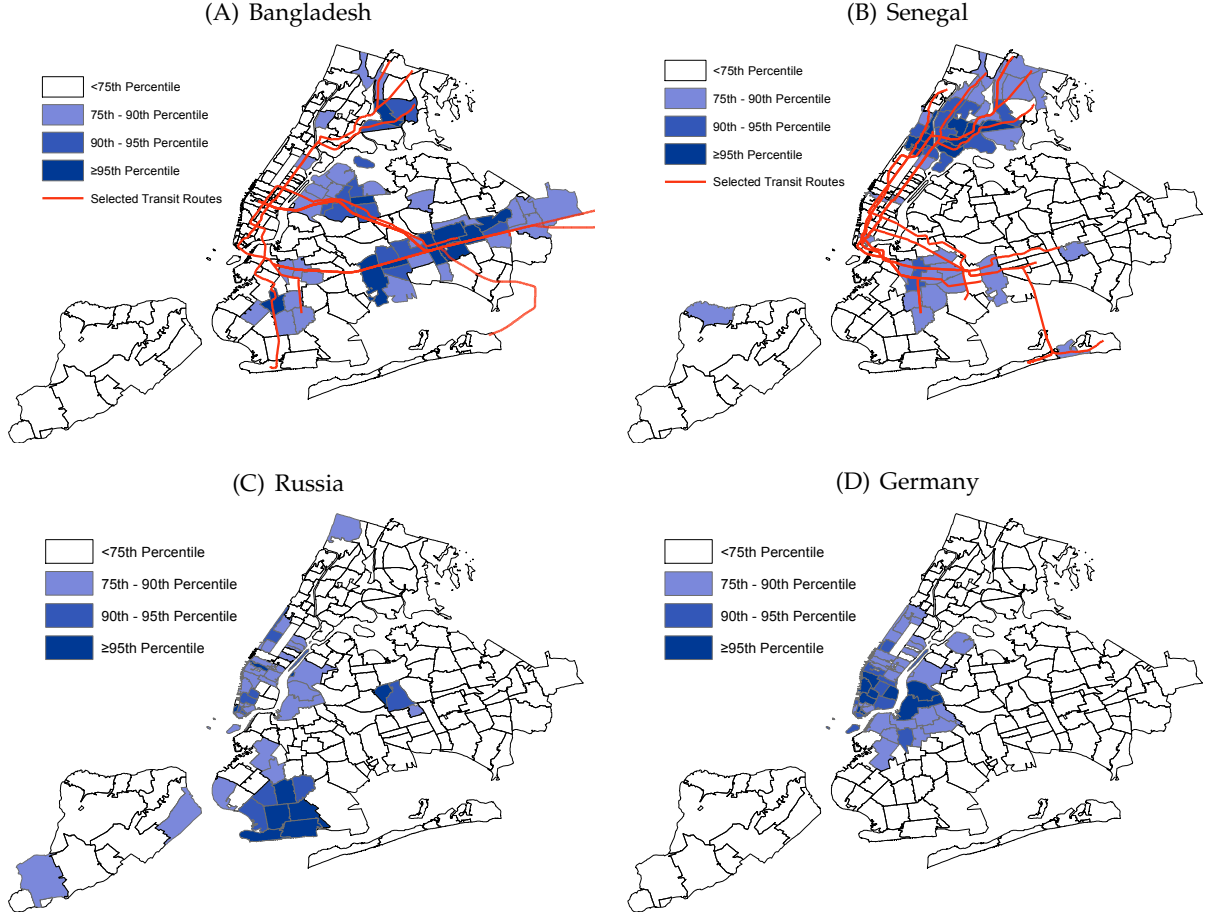
Panel B of Figure 6 shows the resultant connected communities. At the CSA level, the algorithm groups the majority of Long Island with NYC. New Jersey’s border with Pennsylvania and New York is mostly preserved, with the exception of a patch of New York that is grouped with northern New Jersey (in this area, several New Jersey Transit and MTA Metro North lines in New Jersey extend north into New York). Connecticut’s border is also largely preserved. Both upstate New York and Pennsylvania are broken into numerous smaller communities. Unlike the connected communities constructed within NYC, all connected communities at the CSA level are contiguous.

## 5 International Dimension of Social Networks

In addition to exploring the domestic social connectedness of the New York metro area, we next look at the international dimension of social networks in New York. Figure 7 shows the percentile rank of the probability that a user in a given zip code within NYC has of being connected on Facebook to a user in a given country. Panel A shows connections to Bangladesh. Those areas with a high degree of social connectedness to Bangladesh correspond to areas within NYC with large Bangladeshi populations (NYU School of Medicine, NYU Center for the Study of Asian American Health, 2019). Notably, regions with strong connectivity to Bangladesh are almost entirely concentrated along the LIRR and the 4-5-6 train service, consistent with a desire among recent immigrant communities to live in areas with easy travel to existing ethnic enclaves. Panel B shows connections to Senegal and reveals a distribution of connections concentrated in Harlem and portions of Brooklyn north of Prospect Park. Both of the locations contain a substantial number of residents with Senegalese backgrounds (Duthiers, Chen, and CNN, 2013; All Peoples Initiative, 2009). There are similarities to the distribution of the

social network of East Harlem zip code 10035 shown in Panel B of Figure 1. This suggests that areas that have strong connections to the same foreign countries are also more likely to be connected with each other. Consistent with this interpretation, the two broad areas with strong connections to Senegal were also grouped together into a joint “connected community” in Figure 6.

**Figure 7: International Connectivity, NYC**

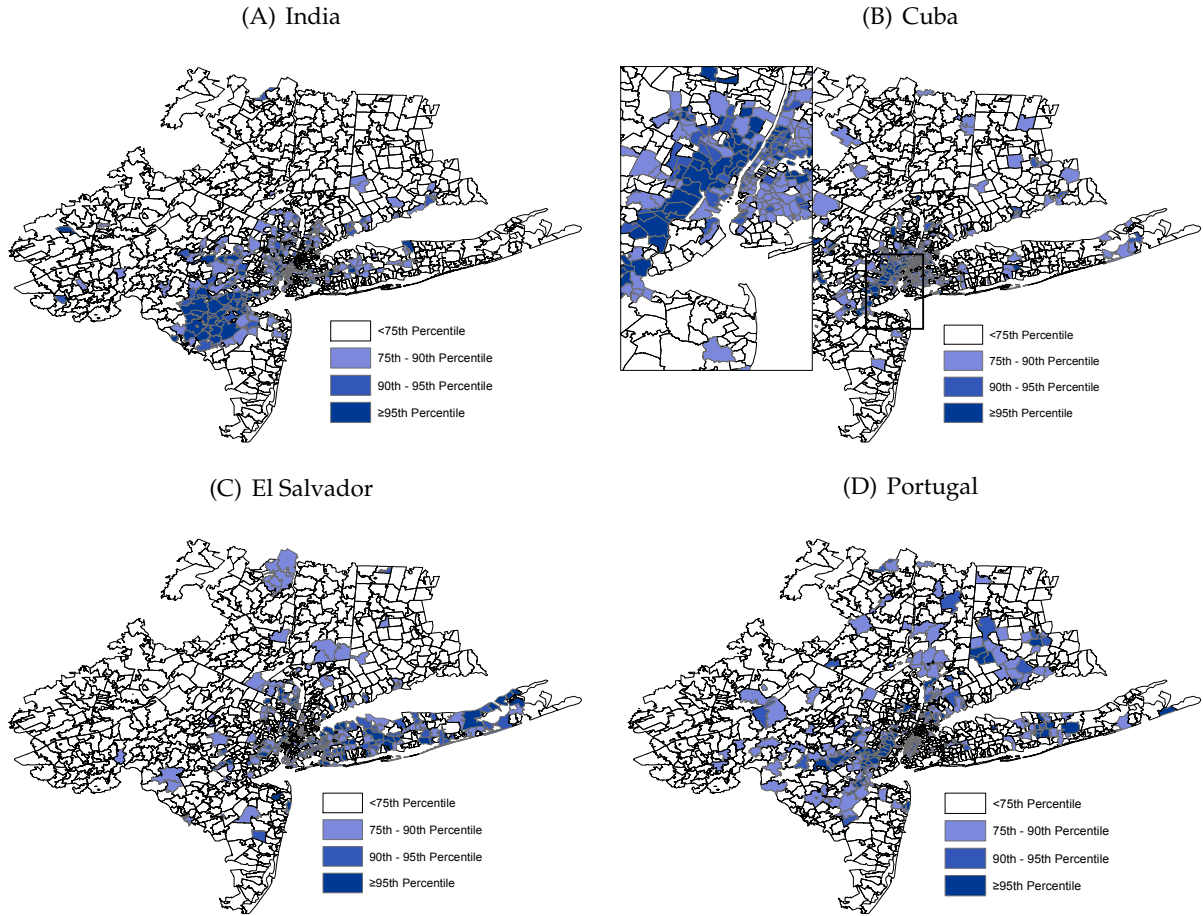


**Note:** Figure shows the percentile rank of the probability of a friendship link, as measured by  $SocialConnectedness_{i,j}$ , of all zip codes  $j$  in NYC to four countries  $i$ : Bangladesh (Panel A), Senegal (Panel B), Russia (Panel C), and Germany (Panel D).

Panels C and D of Figure 7 present the social connectedness of NYC zip codes to two European countries. As shown in Panel C, connections to Russia are concentrated in south Brooklyn, particularly around Coney Island and Brighton Beach; these areas correspond to parts of the city that welcomed large numbers of Russian-speaking immigrants since the 1970s, with increasing numbers arriving after the breakup of the Soviet Union (Ortiz and Untapped Cities, 2014). This pattern of connections to Russia is mirrored by the distribution of connections to most other Eastern European and former Soviet countries. In comparison, the distribution of connections to Germany, which is shown in Panel D, is typical of most Western European countries. Connections to Germany are primarily concentrated in the Midtown and Downtown regions of Manhattan and in neighboring areas in

Brooklyn.

**Figure 8:** International Connectivity, New York CSA



**Note:** Figure shows show the percentile rank of the relative probability of connection, as measured by  $SocialConnectedness_{i,j}$ , of all zip codes  $j$  in the New York CSA to four countries  $i$ : India (Panel A), Cuba (Panel B), El Salvador (Panel C), and Portugal (Panel D).

Figure 8 shows the percentile rank of the probability that a user in a given zip code within the New York CSA has of being connected on Facebook to a user in a given country. Panel A shows connections to India, highlighting the large Indian community in New Jersey (Berger and New York Times, 2008; Batalova, Zong, and New York Times, 2015). Panel B shows connections to Cuba, highlighting the stretch of Cuban communities in New Jersey nicknamed “Havana on the Hudson” (ShareAmerica, 2015). Panel C shows connections to El Salvador, which are primarily concentrated on Long Island. Indeed, El Salvador is the only country with a consulate on Long Island, located in the town of Brentwood (de Relaciones Exteriores de El Salvador, 2019). Panel D shows connections to Portugal revealing several cities and towns referred to as “little Portugal”: Newark, NJ, has the highest concentration of connections to Portugal (Levy and New York Times, 1995), and in New York there are two longstanding immigrant communities on Long Island, Mineola and Farmingville, that display

high degrees of social connectedness to Portugal. Portugal also exhibits high levels of social connectedness to the wealthy northern suburbs, potentially related to vacation travel (Rosenblum and New York Times, 1989; Fishler and New York Times, 2001).

Overall, these findings highlight that the degree of social connectedness of different NYC or New York CSA zip codes is to a substantial degree determined by the presence of migrants from these countries in the respective zip codes.

## 6 Conclusion

We use anonymized and aggregated data from Facebook to better understand the social connectedness of the New York metro area, both at the city level and at the CSA level. We provide evidence for an important role of public transit infrastructure in forming and maintaining urban social connectedness by showing that social networks are distributed along public transportation routes and that social connectedness between locations declines more in travel time than it does in physical distance. We then document a substantial heterogeneity in the geographic concentration of social networks, and highlight that locations with better public transit access have less geographically concentrated social networks, even after controlling for demographic characteristics of the neighborhoods (areas with more geographically dispersed social networks are home to richer and better educated populations). We also show that similarity on socioeconomic characteristics and past migration movements are important drivers of the social connectedness of the New York metro area.

## References

- All Peoples Initiative. 2009. “Senegalese in the New York metro area.”
- Bailey, Michael, Rachel Cao, Theresa Kuchler, Johannes Stroebel, and Arlene Wong. 2018a. “Social connectedness: Measurement, determinants, and effects.” *Journal of Economic Perspectives* 32 (3):259–80.
- Bailey, Michael, Ruiqing Cao, Theresa Kuchler, and Johannes Stroebel. 2018b. “The economic effects of social networks: Evidence from the housing market.” *Journal of Political Economy* 126 (6):2224–2276.
- . 2018c. “The economic effects of social networks: Evidence from the housing market.” *Journal of Political Economy* 126 (6):2224–2276.
- Bailey, Michael, Eduardo Dávila, Theresa Kuchler, and Johannes Stroebel. 2017. “House price beliefs and mortgage leverage choice.” Tech. rep., National Bureau of Economic Research.
- Bailey, Michael, Drew Johnston, Theresa Kuchler, Johannes Stroebel, and Arlene Wong. 2019. “Peer Effects in Product Adoption.” *Working Paper*.
- Bairoch, Paul. 1991. *Cities and Economic Development: From the Dawn of History to the Present*. The University of Chicago Press.
- Batalova, Jeanne, Jie Zong, and New York Times. 2015. “Indian Immigrants in the United States.”

- Baum-Snow, Nathaniel. 2013. "Urban Transport Expansions, Employment Decentralization, and the Spatial Scope of Agglomeration Economies." *Working Paper* .
- Berger, Joseph and New York Times. 2008. "A Place Where Indians, Now New Jerseyans, Thrive."
- Bramouille, Yann, Andrea Galeotti, and Brian Rogers. 2016. *The Oxford Handbook of the Economics of Networks*. Oxford University Press.
- Brooks, Leah and Bryon Lutz. 2014. "Vestiges of Transit: Urban Persistence at a Micro Scale." *Working Paper* .
- Büchel, Konstantin and Maximilian von Ehrlich. 2016. "Cities and the structure of social interactions: Evidence from mobile phone data." *Working Paper* .
- Chin, Seungwoo, Matthew E. Kan, and Hyungsik Roger Moon. 2017. "Estimating the Gains from New Rail Transit Investment: A Machine Learning Tree Approach." *Working Paper* .
- Davis, Donald R, Jonathan I Dingel, Joan Monras, and Eduardo Morales. 2017. "How Segregated is Urban Consumption?" .
- de Relaciones Exteriores de El Salvador, Ministerio. 2019. "Consulado General de El Salvador en Long Island, New York."
- Duggan, Maeve, Nicole B Ellison, Cliff Lampe, Amanda Lenhart, and Mary Madden. 2015. "Social media update 2014. Pew Research Center."
- Duggan, Maeve, Shannon Greenwood, and Andrew Perrin. 2016. "Social media update 2016. Pew Research Center."
- Duthiers, Vladimir, Adeline Chen, and CNN. 2013. "Little Senegal in the Big Apple: Harlem's West African heart."
- Fishler, Marcelle S. and New York Times. 2001. "Long Island Journal: Long Islanders, Yes, But Not Quite Home."
- Glaeser, Edward. 2005. "Urban Colossus: Why is New York America's Largest City?" *Working Paper* .
- . 2011. *Triumph of the City*. Pan.
- Glaeser, Edward and Joshua Gottlieb. 2009. "The Wealth of Cities: Agglomeration Economies and Spatial Equilibrium in the United States." *Journal of Economic Literature* 47 (4):983–1028.
- Glaeser, Edward, Hyunjin Kim, and Michael Luca. 2017. "Nowcasting the Local Economy: Using Yelp Data to Measure Economic Activity." *Working Paper* .
- Glaeser, Edward and Jesse Shapiro. 2001. "Is There a New Urbanism? The Growth of U.S. Cities in the 1990s."
- Glaeser, Edward and Bryce Millett Steinberg. 2016. "Transforming Cities: Does Urbanization Promote Democratic Change?" *Working Paper* .
- Glaeser, Edward L and Matthew E Kahn. 2004. "Sprawl and urban growth." In *Handbook of regional and urban economics*, vol. 4. Elsevier, 2481–2527.

- Glaeser, Edward L., Hedi D. Kallal, José A. Scheinkman, and Andrei Shleifer. 1992. "Growth in Cities." *Journal of Political Economy* 100 (6):1126–1152.
- Granovetter, Mark. 2005. "The impact of social structure on economic outcomes." *The Journal of Economic Perspectives* 19 (1):33–50.
- Granovetter, Mark S. 1977. "The strength of weak ties." In *Social networks*. Elsevier, 347–367.
- Herrera-Yague, C, C.M. Schneider, T. Couronne, Z. Smoreda, R.M. Benito, P.J. Zurfiria, and M.C. Gonzalez. 2015. "The anatomy of urban social networks and its implications in the searchability problem." *Scientific Reports* 5.
- Holahan, CJ, BL Wilcox, MA Burnam, and RE Culler. 1978. "Social satisfaction and friendship formation as a function of floor level in high-rise student housing."
- Ioannides, Yannis M. 2015. "Neighborhoods to nations via social interactions." *Economic Modelling* 48:5–15.
- Ioannides, Yannis Menelaos. 2013. *From neighborhoods to nations: The economics of social interactions*. Princeton University Press.
- Jackson, Matthew O. 2014. "Networks in the understanding of economic behaviors." *The Journal of Economic Perspectives* 28 (4):3–22.
- Jacobs, Jane. 1969. *The Economy of Cities*. Vintage.
- Kowald, Matthais, Pauline van den Derg, Andreas Frei, Juan-Antonia Carrasco, Theo Arentze, Kay Axhausen, Diana Mok, Harry Timmermans, and Barry Wellman. 2013. "Distance patterns of personal networks in four countries: a comparative study." *Journal of Transport Geography* 31:236–248.
- Lazarsfeld, P. and R. K. Merton. 1954. "Friendship as a social Process: A Substantive and Methodological Analysis." In *Freedom and Control in Modern Society*, edited by M. Berger, T. Abel, and C. H. Page. New York: Van Nostrand, 18–66.
- Levy, Clifford J. and New York Times. 1995. "A Portuguese Village in Newark."
- Marmaros, David and Bruce Sacerdote. 2006. "How do friendships form?" *The Quarterly Journal of Economics* :79–119.
- MTA. 2016a. "Annual Bus Ridership by Route."
- . 2016b. "Annual Subway Ridership."
- NYU School of Medicine, NYU Center for the Study of Asian American Health. 2019. "The Bangladeshi Community in the United States and New York City."
- Ortiz, Brennan and Untapped Cities. 2014. "NYC's Micro Neighborhoods: Little Odessa in Brighton Beach, Brooklyn."
- Perlman, Elizabeth Ruth. 2016. "Dense Enough To Be Brilliant: Patents, Urbanization, and Transportation in Nineteenth Century America." *Working Paper* .
- Picard, Pierre M and Yves Zenou. 2018. "Urban spatial structure, employment and social ties." *Journal of Urban Economics* 104:77–93.
- Rosenblum, Ken and New York Times. 1989. "Portuguese Voyagers Reach New Shores."

- Schläpfer, Markus, Luís M. A. Bettencourt, Sébastien Grauwin, Mathias Raschke, Rob Claxton, Zbigniew Smoreda, Geoffrey B. West, and Carlo Ratti. 2014. "The scaling of human interactions with city size." *Journal of The Royal Society Interface* 11 (98).
- ShareAmerica. 2015. "Havana on the Hudson: How Cuban Americans remade Union City, New Jersey."
- Taxi, NYC and Limousine Commission. 2016. "TLC Trip Record Data."
- Verbrugge, Lois M. 1983. "A Research Note on Adult Friendship Contact: A Dyadic Perspective." *Social Forces* 62 (1):78–83.
- Zipf, George Kingsley. 1949. *Human behavior and the principle of least effort*. addison-wesley press.