

1 **Measuring the impacts of economic well being in**
2 **commuting networks — A case study of Bogota,**
3 **Colombia**

4 Manuel A. Florez
5 Department of Earth, Atmospheric and Planetary Sciences, MIT
6 77 Mass. Ave., Cambridge, MA 02139
7 Tel: +1(617) 3088716; Email: mflorez@mit.edu

8 Shan Jiang
9 Department of Civil & Environmental Engineering, MIT
10 77 Mass. Ave., Cambridge, MA 02139
11 Tel: +1(857) 654-5066; Email: shanjiang@mit.edu

12 Ruiqi Li
13 School of Systems Science, Beijing Normal University
14 No. 19 Xijiekouwai St., Haidian District, Beijing 100875, China
15 Department of Civil & Environmental Engineering, MIT
16 77 Mass. Ave., Cambridge, MA 02139
17 Tel: +1(857) 891-3655; Email: liruiqi@mit.edu

18 Carlos H. Mojica
19 Head of Strategic Planning, Service Planning and Technical Affairs
20 Transmilenio S.A., Bogota, Colombia
21 Email: carlos.mojica@transmilenio.gov.co

22 Ramiro A. Rios
23 Inter-American Investment Corporation
24 1300 New York Avenue, N.W. Washington, D.C. 20577, USA
25 Tel: +1(202) 623-2399; Email: rarios@iadb.org

26 Marta C. González *
27 Department of Civil & Environmental Engineering, MIT
28 77 Mass. Ave., Cambridge, MA 02139
29 Tel: +1(857) 928-4546; Email: martag@mit.edu

30 * Corresponding author

31 Word count: 5133 words + 8 figures × 250 words + 1 table × 250 words = 7383

32 Submitted: November 16, 2016

1 **ABSTRACT**

2 Big data such as call detail records (CDRs) from mobile phones are novel resources for travel de-
3 mand models. An important open question is how to use them to extract practical information in
4 relation to urban mobility, socioeconomic development, and well-being. Can we study individual
5 mobility characteristics by income group through the lens of Big Data? In this paper, we present
6 a data analysis framework that uses urban mobility extracted from CDRs, to study various charac-
7 teristics of the commuting network of Bogota, Colombia, relating them to income groups by their
8 residential location. We show that the diversity of commuting trips, defined in terms of entropy of
9 the trips, increases with the income of the population. Further, we show that vehicle travel times
10 during commuting hours from lower income groups clearly suffer longer congested travel times.
11 Our results detail a method to use passively generated mobile phone data as a low cost alternative
12 for transportation policies that can benefit from economic well-being measures for population with
13 different income levels.

1 INTRODUCTION

2 Rapid urbanization has become a common theme across the world, and its influence is profound
3 in Latin America, where over 80% of the population lives in cities. Massive urban migration has
4 imposed enormous burdens to the existing infrastructure, leading to increased congestion, longer
5 travel time and delays, and severe environmental degradation. Spatial segregation is also a salient
6 issue (1). Oftentimes, unintended effects of government policies have further accentuated this
7 problem. For example, stratification became widely spread in major Colombian cities, and it was
8 in fact sectioned as law in 1994 (nueva ley de servicios publicos). Of particular interest is the case of
9 Bogota, the capital and largest city of Colombia, where careful socioeconomic stratification—the
10 basis of a cross subsidization scheme of public utilities for low income residents (2, 3)— has been
11 implemented since the Eighties. New insights into the impact of this kind of policy on individual
12 mobility and socioeconomic well-being can be obtained through the lens of Big Data.

13 Latin American cities, as those in much of the developing world, exhibit sharp differences
14 in income levels among various neighborhoods and are characterized by economic and urban spa-
15 tial segregation, closely related to job opportunities, economic mobility, and travel behavior (1, 2).
16 There has been a lot of qualitative work on understanding spatial segregation (3, 4, 5, 6, 7), how-
17 ever, there are still ample opportunities to better quantify it using new data resources. In particular,
18 it is necessary to come with effective alternatives at low cost to measure impacts of the spatial
19 organization of cities on the well-being of their inhabitants. Understanding the interplay between
20 economic segregation and mobility as extracted from communication technologies is important for
21 more efficient means of transportation and urban planning (8, 9).

22 To that end, we propose a framework of analysis to use big data through the lens of urban
23 transportation (10) to quantify the impact of economic segregation on individual mobility. We
24 utilize anonymized call detail records (CDRs) with an observation period of 6 months (in 2013
25 to 2014) to estimate the origin-destination (OD) matrices for a typical day in the city of Bogota,
26 Colombia. We validate the ODs against the estimates from local official household travel survey
27 (2011). We then create sampled networks with around 20,000 travelers for each of the 4 income
28 brackets, over which we estimate the duration of home-based-work (HBW) trips by income group
29 (socioeconomic stratum). We show differences in their mobility diversity and congested travel
30 times. We demonstrate that the impact of segregation on urban mobility can be extracted from
31 CDR data, which can be used as a good alternative to time consuming and expensive travel sur-
32 veys.

33

34 LITERATURE REVIEW

35 Spatial segregation in its most extreme form originated in the ghetto, a clustering in space of iden-
36 tifiable ethnic groups, often times instituted by political authorities. The causes and impacts of
37 modern spatial segregation on society are a matter of debate (3, 4, 5, 6, 7) . Vandell et al. (11) ar-
38 gue that in the urban space, segregation is the consequence of at least three factors: administrative
39 policy, market forces, and individual choices. While many studies have shown the socioeconomic
40 consequences of segregation, few have focused on using information technologies to measure the
41 impacts of segregation on urban mobility. The available information has been limited (12), restrict-
42 ing the sizes of representative samples by income groups, particularly for residents in poverty.

43 In order to improve the mobility of a city, transportation and urban planners need to quantify
44 the interplay between travel demand and existing infrastructures. This is typically done via models

1 of travel demand, which estimate daily trips of individuals aggregated to origin-destination (OD)
2 matrices (13) by mode, purpose, and time of the day. The seed information to create these trip
3 matrices need to capture travel preferences of various population groups across the day. This
4 requires careful sampling of preferences through representative travel diaries. While these surveys
5 are rich in detail per individual, they are expensive and quickly become outdated. They only
6 contain one or two days for a thin sample (usually 1 percent households in a metropolitan area) to
7 model trips over the years in a city with millions of residents (12, 14). These limitations are very
8 problematic as local municipal budgets are increasingly constrained and people’s mobility patterns
9 are dynamic. As a result, cheaper and more up-to-date data sources and analytical methods are
10 called to increase the efficacy of urban travel demand models to empower urban and transportation
11 planning in the age of information and telecommunications.

12 Recent research has explored the opportunity to utilize cell phone data, known as *call*
13 *detail records* (CDRs) for modeling travel behavior. CDRs are metadata on phone usage (such
14 as phone calls, data or text messaging) collected by a cellular carrier for billing and operational
15 purposes, and contain geospatial whereabouts of users during their phone usage. As cell phones
16 have become ubiquitous, a surge of research (15) have demonstrated that various insights on human
17 behavior can be extracted from the massive, passively collected CDR datasets. Techniques have
18 been used to identify daily mobility motifs, which have simplified the extraction of daily trip chains
19 of individuals without using surveys (16). Mobile phone meta data have also changed the process
20 of modeling the spreading of infectious diseases such as dengue and malaria (17, 18, 19).

21 The main advantage of individual phone records is that they contain hidden valuable in-
22 formation on the most visited locations by each user across time. Mobile market share typically
23 covers a great fraction of the population during various months and each individual tends to visit
24 few repeated locations in their journeys (20). While the geospatial tags are not accurate in space
25 and time to generate complete journeys in a day, depending on the data collection technologies,
26 treated with right methods, it is possible to obtain average transportation demand matrices by pur-
27 pose and by hour-of-the-day that represent the travel demand of an entire city (10, 13, 21).

28 Of particular interest in this domain is the estimation and validation of OD matrices. Early
29 works presented transient OD matrices using CDRs (22, 23) which mapped road usage to the
30 neighborhoods that originated the travels. But their validation with existing models that also in-
31 cluded mode, routes and travel times remained a challenge. In the last two years, the techniques to
32 generate ODs from phone data as well as their validation have been further developed. Using only
33 CDRs and census population data, (21) and (10) demonstrate techniques that can estimate ODs
34 by purpose (e.g. commuting trips, or home-based-other) and by time-of-day (AM, PM etc.). The
35 resulting ODs in Boston and Rio de Janeiro matched the ODs from existing models that required
36 survey data. The work in (13) describes a system architecture for an end-to-end software solution
37 that transforms and integrates mobile phone data into estimates of travel demand and infrastruc-
38 ture performance, and applies it in five cities comparing favorably with existing models based on
39 surveys. CDRs and population data are consistently transformed into OD matrices by purpose and
40 time of day, and routes through road networks are constructed using open and crowd-sourced data
41 repositories. The analytics on the system’s output is fast and portable.

42 In the following sections, we first describe the data used in this study. We then employ
43 similar methods tested for other cities (10, 13, 21) to transform CDRs into ODs and validate the
44 modeling results with local household travel survey data for Bogota, Colombia (12). Next, we
45 present a network analysis method to enable the understanding of the interplay between economic

1 segregation and individual mobility. The main challenge is how to obtain accurate transportation
2 information by aggregating trips without losing the representation by income from origins of the
3 trips. While the technique is applied from end to end in Bogota, the results are portable to any
4 other region with census and CDR data available.

5 **DATA AND METHODOLOGY**

6 **Data**

7 To estimate travel demand for population by income stratum in Bogota, Colombia, we obtained
8 (i) an anonymized CDR dataset from a telecommunication operator in Colombia, including in-
9 formation for 1.5 million users for a period of 6 months across 2013 and 2014, (ii) population at
10 the census block level with their income stratum (ranging from 1 to 6, representing income level
11 from low to high). To validate our estimate results, we also obtained: (iii) a latest set of household
12 travel survey data (in 2011) for transportation planning purpose, officially authorized by the city's
13 department of transportation and conducted by an international transportation consulting firm, and
14 (iv) a set of car travel times for the estimated OD matrices at the tower-level in a weekday morning
15 peak hour (7:00 am -8:00 am), queried from the API of an online mapping service. We discuss the
16 details of these data sets in the following subsections.

17 *Call Detail Records*

18 The CDR dataset was gathered at the cellular tower-level, with 659 towers distributed across the
19 Bogota metropolitan area (with an area of 477 square kilometers, and a population density of
20 16,143 persons per sq. km.). The study area encompassing the Capital District with its 20 locali-
21 ties (localidades) and the neighboring municipalities of Soacha, Mosquera, Funza, Madrid, Chía,
22 Cajicá, Cota, La Calera, Tenjo, Tabio, Sibaté, Zipaquirá and Facatativá, for 912 transportation
23 analysis zones (ZATs). Each record in the CDR dataset contains an anonymous user ID, the ge-
24 ographical location in the form of the latitude and longitude of the cellular tower, and the time at
25 the instance of the phone activity. We analyze the OD trips at both cellular tower level for network
26 analysis purpose and at the locality and ZATs level for OD validation purposes. Figure 1 A-C.
27 shows the basic statistics on the CDR dataset.

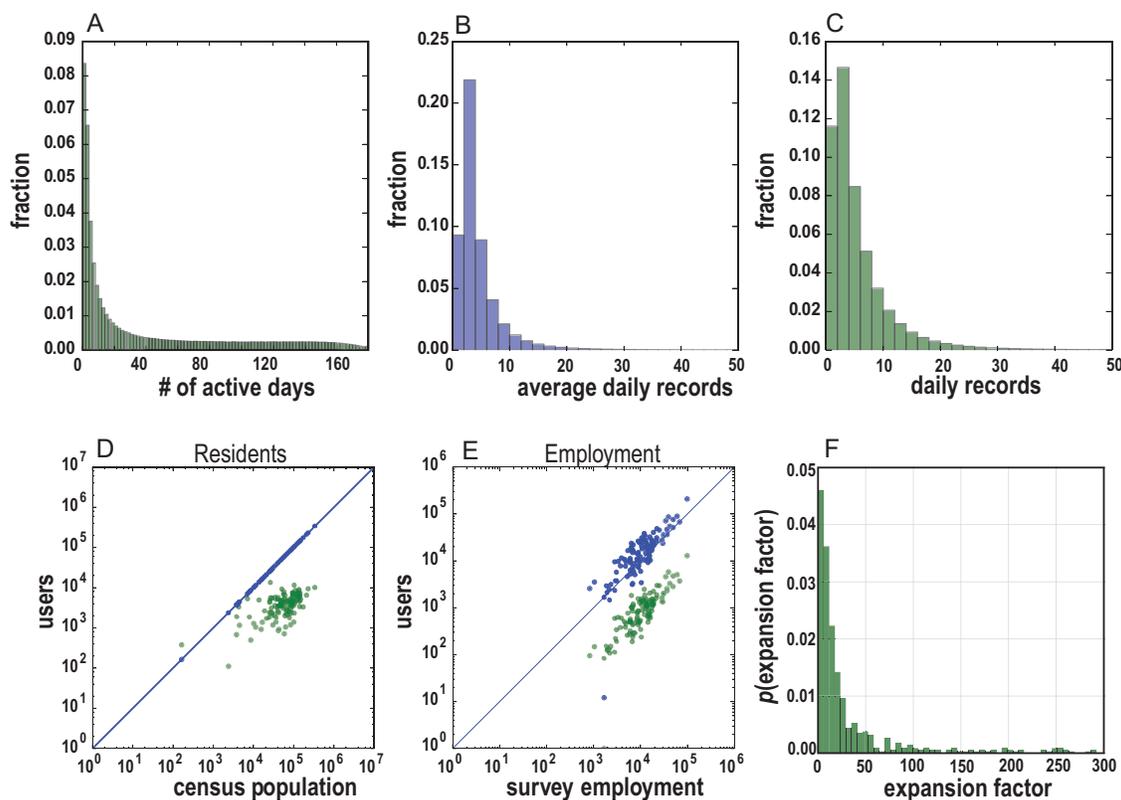


FIGURE 1 A-C. Distribution of number of active days of the Bogota phone data. B. Distribution of the individual users' average daily records. C. Distribution of daily phone records for all phone users over the 6-month observation period. D. Correlation between the WorldPop residential population and the mobile phone estimated residents before (green) and after (blue) expansion adjustment. E. Correlation between the employment population estimated from local household travel survey data and that estimated using the mobile phone data. F. Distribution of the expansion factors which are used to expand mobile phone users to total population.

1 *Population and Income*

2 To understand the economic status of local communities, we obtained a data-set at the city block
 3 level. The total population in Bogota is over 7 million. The census population data were also used
 4 to help expand the CDR sample users to the population for the Bogota metropolitan area, in order
 5 to generate representative urban travel demand estimates. To do so, in Fig. 1 D, we show the before
 6 (in green) and after (in blue) comparison of CDR users and census population at the ZAT level for
 7 Bogota. The purpose is to calculate expansion factors at the geographic level of analysis (e.g.,
 8 shown here at the ZAT level). Fig. 1 F shows the distribution of expansion factors over Bogota.

9 We use the official socioeconomic stratification as a proxy for the spatial distribution of in-
 10 come within the city of Bogota. The department of city planning (DCP) is legally responsible for
 11 assigning a socioeconomic stratum to each city block (2). A scale from 1 through 6 is used, where
 12 1 corresponds to the lowest socioeconomic level and 6 to the highest. The socioeconomic level of
 13 each city block is determined by a DCP official who relies on direct observation of the block and
 14 its surroundings and must take into account the following factors in his assessment (24): physical

1 characterizes of buildings, condition of local roads, presence and quality of sidewalks, ease of ac-
 2 cess to major roads and public transport, quality of urban space surrounding the block and overall
 3 urban context of the neighborhood. Socioeconomic stratification is the basis of a differentiated
 4 pricing scheme for public utilities (2). Residents of blocks classified in the three highest socioeco-
 5 nomic levels pay proportionally higher rates that are used to subsidize residents in the three lowest
 6 levels. Thus, careful assignment and regular updates to the socioeconomic stratification of the city
 7 are guaranteed by the current regulatory framework. Furthermore, in countries such as Colombia,
 8 where the informal economy plays a significant role, income data is grossly under-reported in both
 9 the census and local surveys (25). Using socioeconomic stratification as a proxy for income is a
 10 reasonable alternative that does not suffer from such biases.

11 There are about 45 thousand blocks classified by income (socioeconomic stratification) in
 12 the city but only 659 cell-towers. We need to assign an income rank to each cell-tower coverage
 13 area. This requires careful aggregation of the block-level income categories to the cell-tower scale.
 14 We map each city block to its corresponding cell-tower coverage area by calculating the centroid
 15 of the block and determining the converge area it falls within. Then, for each cell tower coverage
 16 area we calculate the population-weighted income rank:

$$\mu_n = \sum_{i=1}^{l_n} w_i s_i, \quad (1)$$

17 where l_n is the number of blocks that fall within coverage area n , w_i is the population weight of
 18 block i , defined as the fraction of the total population of block i , s_i is the discrete socioeconomic
 19 stratum (from 1 to 6) assigned to block i by the city planning officials.

20 *Survey of Mobility*

21 The 2011 Bogota Survey of Mobility (12, 26) was a local household travel survey, including 45
 22 thousand individual samples to represent the over 7 million Bogota residents. The survey was con-
 23 ducted to collect trip information for residents during one sample day, including their trip purpose,
 24 and departure and arrival time and zonal information, and their social demographic information.
 25 Although there are more than 900 ZAT zones in the Bogota metropolitan area, due to the limited
 26 sample size, the survey only covered residents' travel in 767 ZATs. Using the estimates from the
 27 survey data, we validate our estimated employment population by inferring cellphone users' work
 28 location (discussed later), combined with our calculated expansion factors at the ZAT level (shown
 29 in Fig. 1 F). Fig. 1 E shows high correlation between the employment estimates from the cellphone
 30 data model and those from the model informed by the survey data.

31 *Travel Times in an AM Peak Hour*

32 To assess the impact of spatial segregation on travel times, we estimate the trip duration for OD pairs
 33 that have a high number of home-based work trips. Bogota is a very congested city; estimates of
 34 trip duration have to take into account the increase in travel times during rush hours. To do so, we
 35 use the estimates of congested travel times provided by an online mapping service API (e.g., the
 36 Google traffic API, source: <https://developers.google.com/maps/documentation/directions>).

37 **Generating OD Flows from CDR Data**

38 Combined with population data and the CDR data, we generate ODs by trip purpose and by time-
 39 of-day based on three key steps: stay detection, activity labeling, and trip generation, as discussed

1 in (21).

2 *Stay Detection*

3 In order to discover users' activity location, we first filter out noise resulting from tower-to-tower
4 call balancing performed by the mobile service provider, creating the appearance of false move-
5 ments. We then employ a method (27) to distinguish users' stationary stay locations (when/where
6 users engage in an activity) from their moving pass-by locations (when/where users are en-route
7 to activities). As the data is tower-based, one can only know the closest tower to the user's actual
8 location, so the estimate of a user's position is known up to the Voronoi cell for that tower. Due
9 to the discrete nature of this data, the aforementioned call sequence simplification is carried out
10 by joining sequences of calls made from a sets of towers within a certain distance threshold, fol-
11 lowed by joining the sequence of calls made from the same tower. To address issues of temporal
12 resolution, we only keep stays if the user is known to be in that location for at least 10 minutes.

13 *Activity Labeling*

14 To successfully extract purposes for every trip, we classify activities as *home*, *work* and *other*.
15 Human mobility patterns as captured from mobile phone data that exhibit regularity and frequent
16 returns to previously visited sites. For every user, her most visited location on weekdays from
17 7pm-8am and on weekends is classified as her *home*. Users with too little activity from their home
18 locations are filtered out of the analysis. This is followed by assigning the user's *work* place, which
19 is defined to be the non-home location that the user visits second most during the complement
20 of the home time period on weekdays. Similarly, users with too few calls from their assigned
21 workplace are excluded. Stays made from other locations are all classified under *other*. Once each
22 stay is labeled with an activity purpose, then the resulting trips obtained from stay locations can be
23 assigned with purpose pairs, such as home-based-work (HBW), home-based-other (HBO), or non-
24 home-based (NHB). The ODs obtained in this way are then classified in terms of their purpose
25 pairs. These methods are not necessarily definite solutions for perfectly estimating users' home
26 and work locations, but they are straightforward and may lead in some cases to incorrect labeling
27 of home and work locations. However, with increased spatial and temporal granularities of data
28 and the inclusion of refined GIS information, more sophisticated algorithms can be developed.

29 *Trip Estimation*

30 After the call data have been assigned one of the three activity tags (*home*, *work* or *other*), the
31 next step is to go through the time-ordered stay sequence for every user. Two consecutive calls on
32 weekdays constitute a raw trip if they are not from the same location and are in the same effective
33 day, which spans 3am of the previous day to the 3am of the next. Our method assumes that users
34 typically travel from their home location at the beginning of an effective day and travel back home
35 at the end. Therefore if a user's last call of the day is not from the home location, a raw trip is
36 added to home. Similarly, if a user's first call of the day is made from a location other than home,
37 a trip is added to ensure user's travel from her home. As CDR data is passive and generated when
38 users choose to interact with their phones, one cannot assume that users start their trip at the exact
39 time they make the call. We introduce a departure time estimation to account for the passiveness
40 of CDR data as explained in (21). Fig.2 A-D shows the activity start time and duration by activity
41 type of home, work, other, and all types. Fig. 2 E-H compares the trip departure time distribution
42 by trip purpose of HBW, HBO, NHB and all types between the estimates of survey and the CDR-

- 1 data based model. The results show that in general the CDR based model presents similar temporal
- 2 patterns to the survey, although there are some trade-off between the trips by trip purpose.

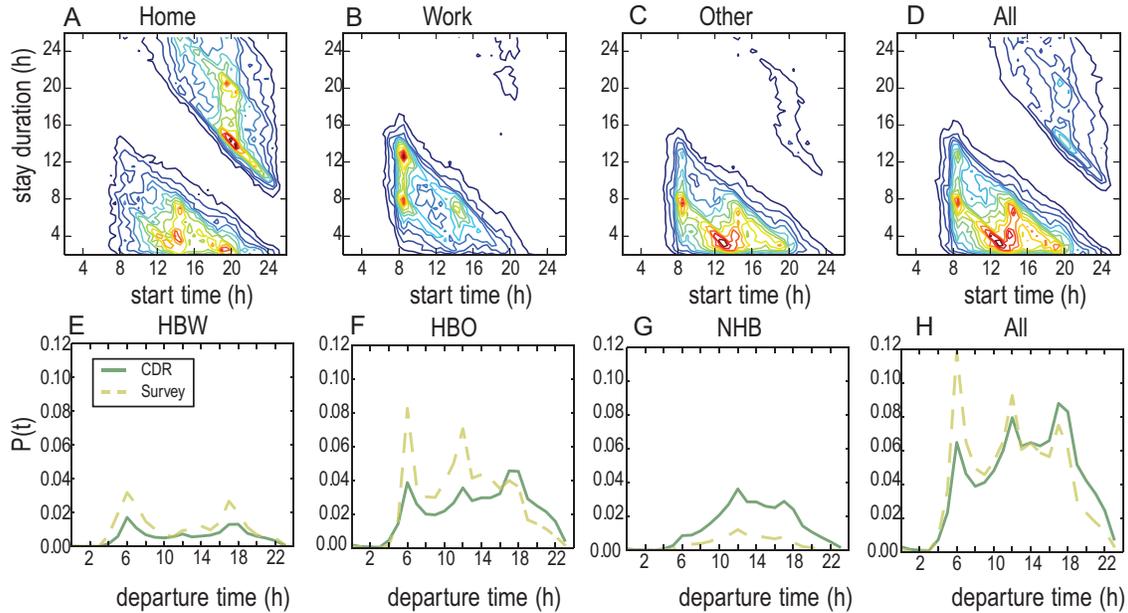


FIGURE 2 A-D. Temporal distribution (start time and stay duration) for inferred activities, including A. Home, B. Work, C. Other, and D. All types of activities. E-H. Comparison of trip start time by trip purpose, estimated with the local household travel survey data (in light green) and with the mobile phone data (in green), for different trip purposes including: E. Home-based-work (HBW), F. Home-based-other (HBO), G. None-home-based (NHB), and H. All trip purposes.

3 *OD Generation*

4 With the estimated trip departure time, and expanding the individual trips, we aggregate trips to
 5 construct OD matrices between zones by counting the number of expanded trips between each pair
 6 of zones on a typical day. We generated ODs for cell-tower based Voronoi-polygons, as well as
 7 for ZATs. We validate the ODs estimated from CDR data against those derived from the mobility
 8 survey at the ZAT level, since the survey only records trip origins and destinations at the ZAT
 9 level. We aggregated OD trips from ZAT to localities in Bogota. Fig. 3 shows the correlation of
 10 the ODs estimated from CDR and from survey data by time of day at the locality level. In general,
 11 the correlations of the CDR-based model and the Survey-based model are very high at the inter-
 12 locality level. The estimates of these two models on the total number of trips per day are reasonably
 13 close (as shown in Table 1). Given the good correlations and the validity of the method for other
 14 cities, we use the CDR-based travel demand model to further analyze the mobility characteristics
 15 by income group, because it is informed by many more users than the survey. This allows us to
 16 sample tens of thousands commuters per income bracket.

17 **Network Analysis**

18 The generated ODs naturally lend themselves to a network structure representation (10, 16, 28).
 19 The weight of the directed link connecting the origin node to the destination node is given by the

Bogota	HBW	NHB+HBO	AM	MD	PM	Total
CDR Trips (millions)	1.53	9.07	1.98	4.10	2.46	10.60
Survey Trips (millions)	2.82	9.80	1.99	5.09	2.48	12.62

TABLE 1 Number of trips on a typical day, split by trip purpose and time of day. AM, MD, and PM refer to 7am to 10am, 10am to 3pm, and 3pm to 7pm.

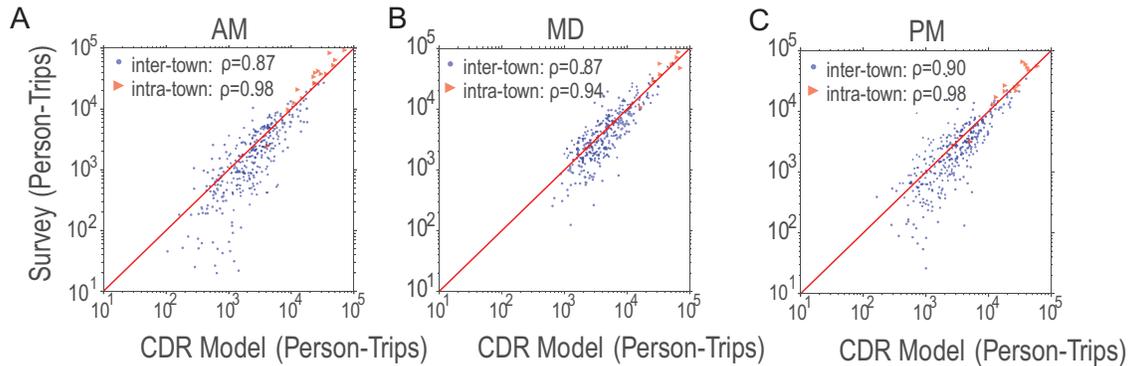


FIGURE 3 OD validation with Household Travel Survey at the Locality Level. by time of day, including A. morning (AM), B. midday (MD), C. evening (PM).

1 number trips estimated for that particular pair. Since we label trips by purpose, we select only
2 OD pairs that represent commutes to work, as pairs where the number of HBW trips is significant.
3 This selection enables us to quantify the relation between economic well being and the structure
4 of the commuting network. The resulting network has 633 nodes and 27,939 links. Figure 5
5 shows the commuting network for Bogota using a geotagged layout, the distribution of degree or
6 number of connections per node (Fig. 5B) and the trips per link (Fig. 5C). We see that while the
7 degree distribution of the origins has an exponential decay, the distribution of the number of trips
8 is broader with few OD pairs having between 100-1000 commuting trips. To assess the impact
9 of income on the different network metrics we divide the nodes into four groups: low income,
10 with weighted income stratum ranging 1.0 to 2.5 contained in 141 nodes, middle income, in the
11 range 2.5-3.5 with 230 nodes, upper middle income, ranging in 3.5-4.5 with 134 nodes and higher
12 income, from 4.5-6.0, and with 128 nodes.

13 In a seminal work Eagle et al. (29) analyzed landline calls and a nationwide mobile phone
14 dataset in the UK. They proposed metrics of topological and spatial network diversity that dis-
15 played a strong correlation with the socioeconomic outcomes of the regional communities that
16 each node represented. In this work we estimate mobility diversity (30) in similar way, first cal-
17 culating the Shannon entropy of each node: $H_i = -\sum_{j=1}^k p_{ij} \log(p_{ij})$, where k is the number of
18 destinations with origin in i , or the degree of node i . p_{ij} is the relative proportion of trips between
19 i and j : $p_{ij} = \frac{T_{ij}}{\sum_{j=1}^k T_{ij}}$. The spatial diversity of node i is then defined as:

$$D_i = \frac{H_i}{\log(k_i)}, \quad (2)$$

20 which is the ratio of the entropy observed in the trips divided by $H^{rand} = \log(k_i)$, the Shannon
21 entropy when all trips from i are weighted equally among all the destinations, meaning $p_{ij} = 1/k_i$.

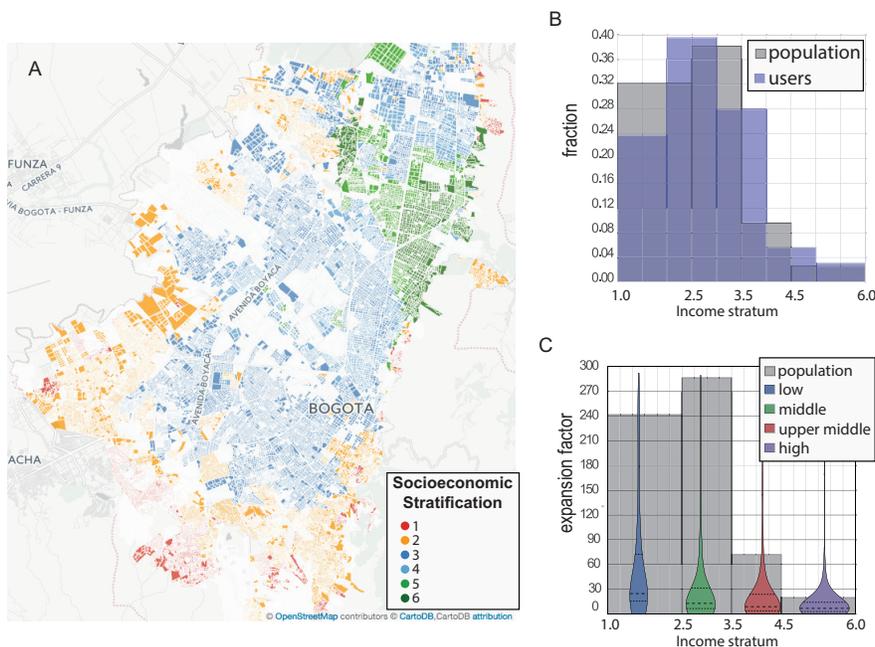


FIGURE 4 A. Map of population distribution by stratum from 1 to 6 going from lowest (1) to highest income (6). B. Population distribution of 4 income brackets (in gray) and fraction of users in each of 6 income strata. Income brackets are defined based on the average stratum on a given trip source. These are low income [1, 2.5], middle income [2.5-3.5], upper middle income [3.5-4.5] and higher income [4.5-6.0]. C. Distribution of factors to expand mobile phone users to population in the different origins separated in four income brackets– in each income bracket we have 1 user each 20 people, showing that the used mobile phone data cover all income levels with little bias.

1 Lower diversity implies that the commuting trips of origin i are more attracted to some destinations
 2 than others, among the k_i observed destinations. Entropy close to 1 implies that the trips are
 3 uniformly balanced. Also, the trip duration and the trip distance per origin are weighted over each
 4 node as follows: $\langle x \rangle_i = \sum_{j=1}^k p_{ij} x_{ij}$, where x_{ij} is the quantity of interest as measured for OD pair
 5 i, j .

6 **RESULTS**

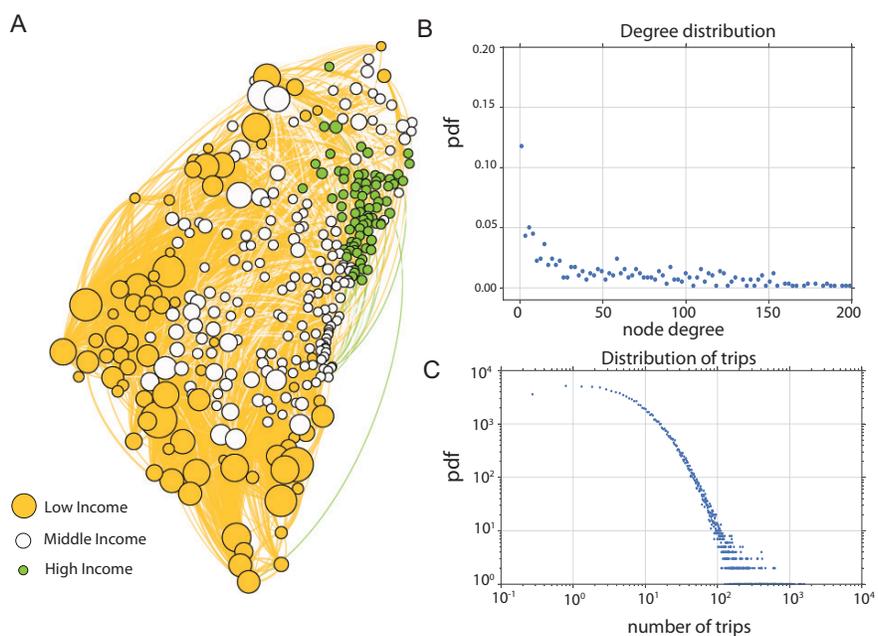


FIGURE 5 Network of commuting trips. **A.** 633 origin nodes colored by income group (low, middle, high) and sizes representing the total number of commuting trips. **B.** Degree distribution represents the number of destinations of each node. **C.** Distribution of trips between the links.

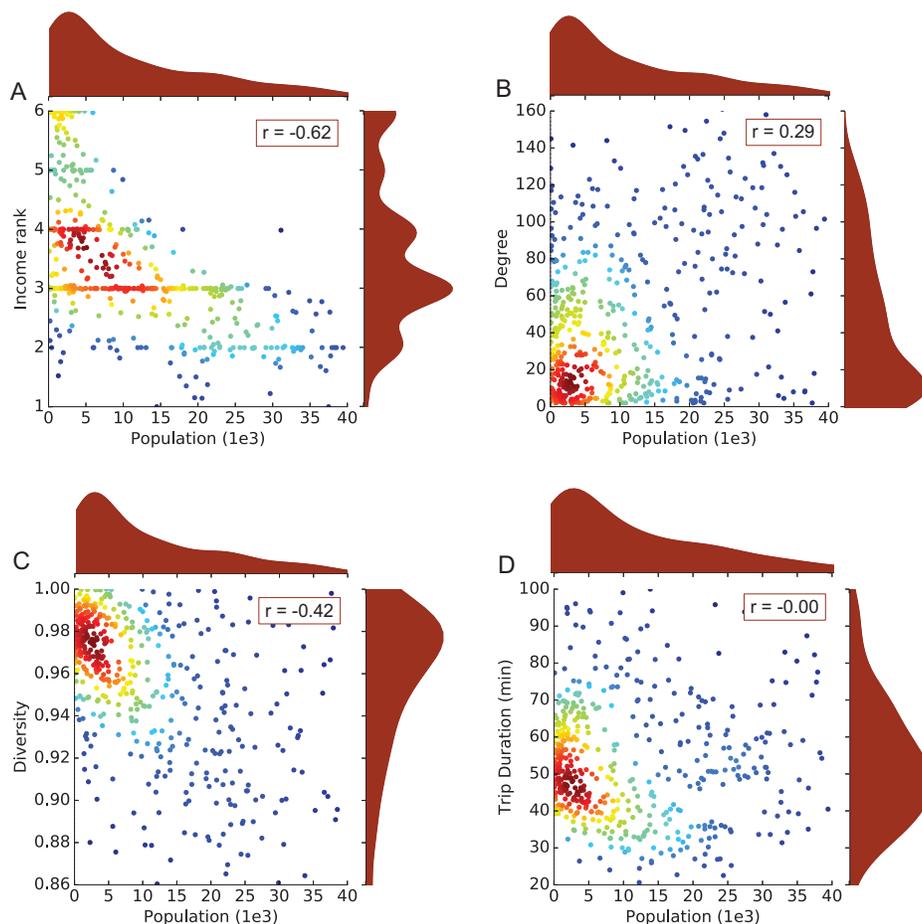


FIGURE 6 Income and various network metrics vs. population of the origins. **A.** Average income **B.** Degree **C.** Mobility diversity **D.** Commuting travel times. Most trip sources have around 5,000 people, those with larger population have lower income and lower mobility diversity. There is not correlation between the average commuting travel time and the population, and surprisingly, very low correlation between the degree and the population of origin.

1 We start by exploring the properties of the network in relation with the population of the
2 origins (Figure 6). Most origins have similar population around 5,000 people. More populated
3 nodes have lower income (Figure 6A) and slightly higher degree (Figure 6B). Origins with larger
4 population have lower diversity (Figure 6C), indicating that they are preferentially attracted to
5 particular destinations. There is not clear relation between the population of the origins and their
6 average travel times in commuting (Figure 6D).

7 Next, we compare diversity and the entropy by income brackets (Figures 7 A-D). We ob-
8 serve that the higher the income the smaller the tendency of the origins to have preferential des-
9 tinations for commuting (Figures 7 A-B), resulting in higher mobility diversity. Indeed, there is
10 significant positive correlation between diversity and income rank (Pearson's r is 0.47). Entropy
11 on the other hand does not display any significant relation with income (Figures 7C-D). In order to
12 obtain estimates independent of the number of travelers per origin, we sample directly the commut-
13 ing trips of 20,000 mobile phone users for each of the 4 income brackets. They are selected with
14 home locations following the population distribution of the city. Similar to the total network, the
15 sampled network has 633 nodes and 25,073 links. The observed effects of entropy and diversity by
16 income group remain for both the sampled network and the entire network, as shown respectively
17 in Figs 7 A vs. B and C vs. D.

18 Finally, we use the sampled network to study the distance to work, trips duration and time
19 spent in congestion by income group (Figure 8). The sampling accounts for the wide differences in
20 population among origins (Figure 6A), leading to similar statistics of the quantities of interest by
21 income group and thus to more meaningful results. For each OD pair, congested travel times and
22 road distances are queried from an online mapping service as discussed in the methods section.
23 The distribution of values for the lowest income bracket tend to be broader. It is striking to see,
24 however, that consistently, people in the lowest income group travel longer distances to work,
25 spend more time in their commutes and are more affected by congestion. In fact, low income
26 residents spend almost twice as much time in congestion than the high income sample (Figure
27 8C).

28 CONCLUSIONS

29 We use mobile phone data and query on-line maps to study the commuting network in Bogota,
30 focusing on the relation between socioeconomic characteristics of the origins and their mobility
31 characteristics. We find that mobility diversity has a significant positive correlation with income
32 rank. Previous studies have confirmed this observation using information and communication
33 technologies (ICTs) in the context of social contacts (29) and trips at a country scale (30). Eagle et
34 al. (29) suggested that the diversity of social contacts is relevant because it showed very high degree
35 of correlation (greater than 0.72) with socioeconomic indicators of well being. More recently,
36 Pappalardo et al. (30) showed that social diversity may not be as important in the predictability of
37 the Gini coefficient of a Municipality, and that the mobility diversity adds the largest predicting
38 power to their regression models. The debate is only starting and many questions remain open,
39 specially for both mobility and social networks within the urban scale, where less studies of this
40 type have been done.

41 To our knowledge, this is the first study showing that mobility diversity is correlated with
42 a socioeconomic indicator of well-being at the lowest possible degree of granularity, the mobile
43 phone tower level. This opens up the possibility of using network based structural metrics from
44 ICTs in the context of urban and transportation planning in relation to economic well being. In

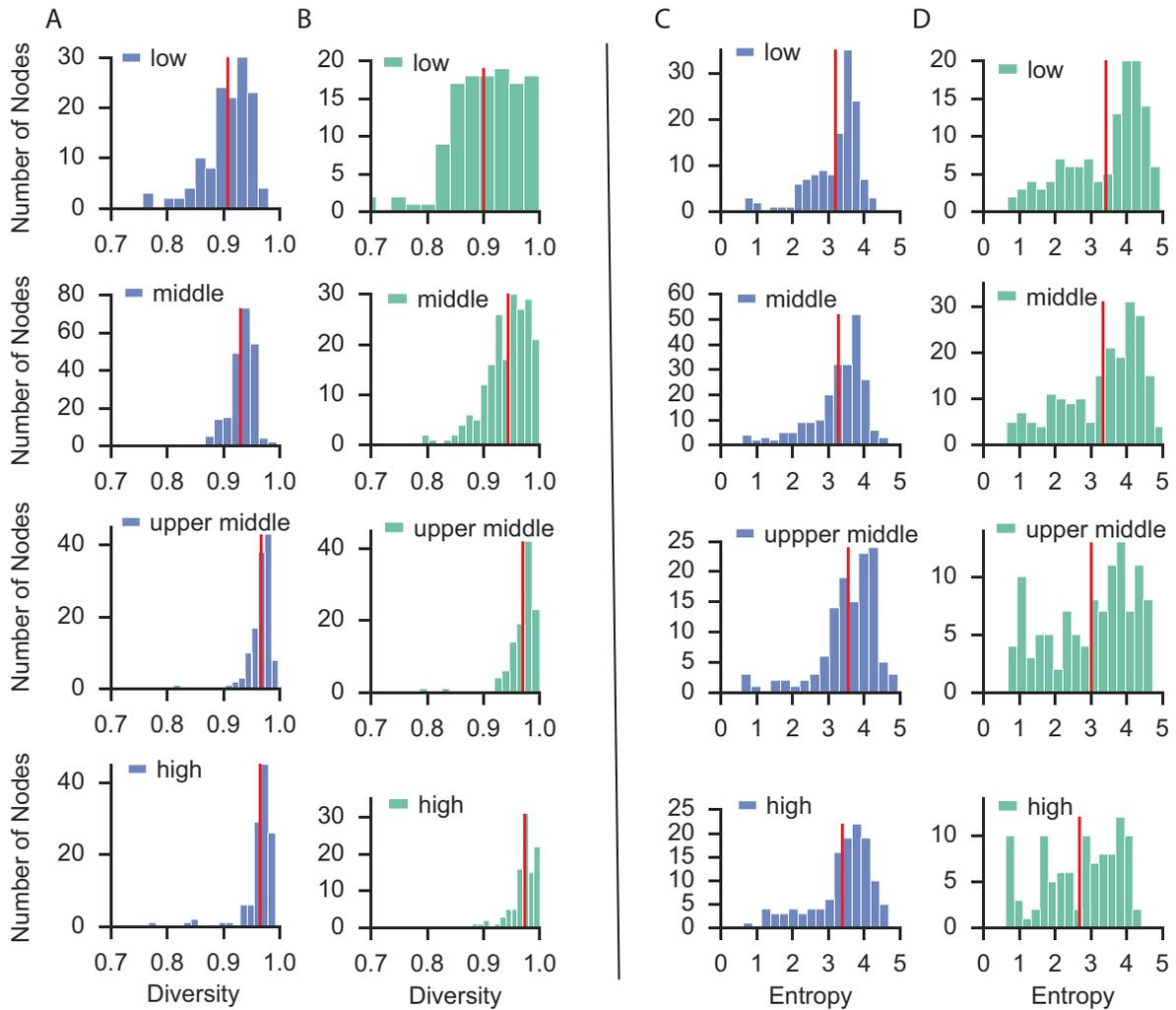


FIGURE 7 Distribution of mobility diversity and entropy by income group. **A.** Diversity in sampled network using same number of travelers per income group **B.** Diversity estimated in the complete commuting network of the CDR-based model **C.** Entropy in sampled network with same number of travelers per income group **D.** Entropy in the complete commuting network of the CDR-based model

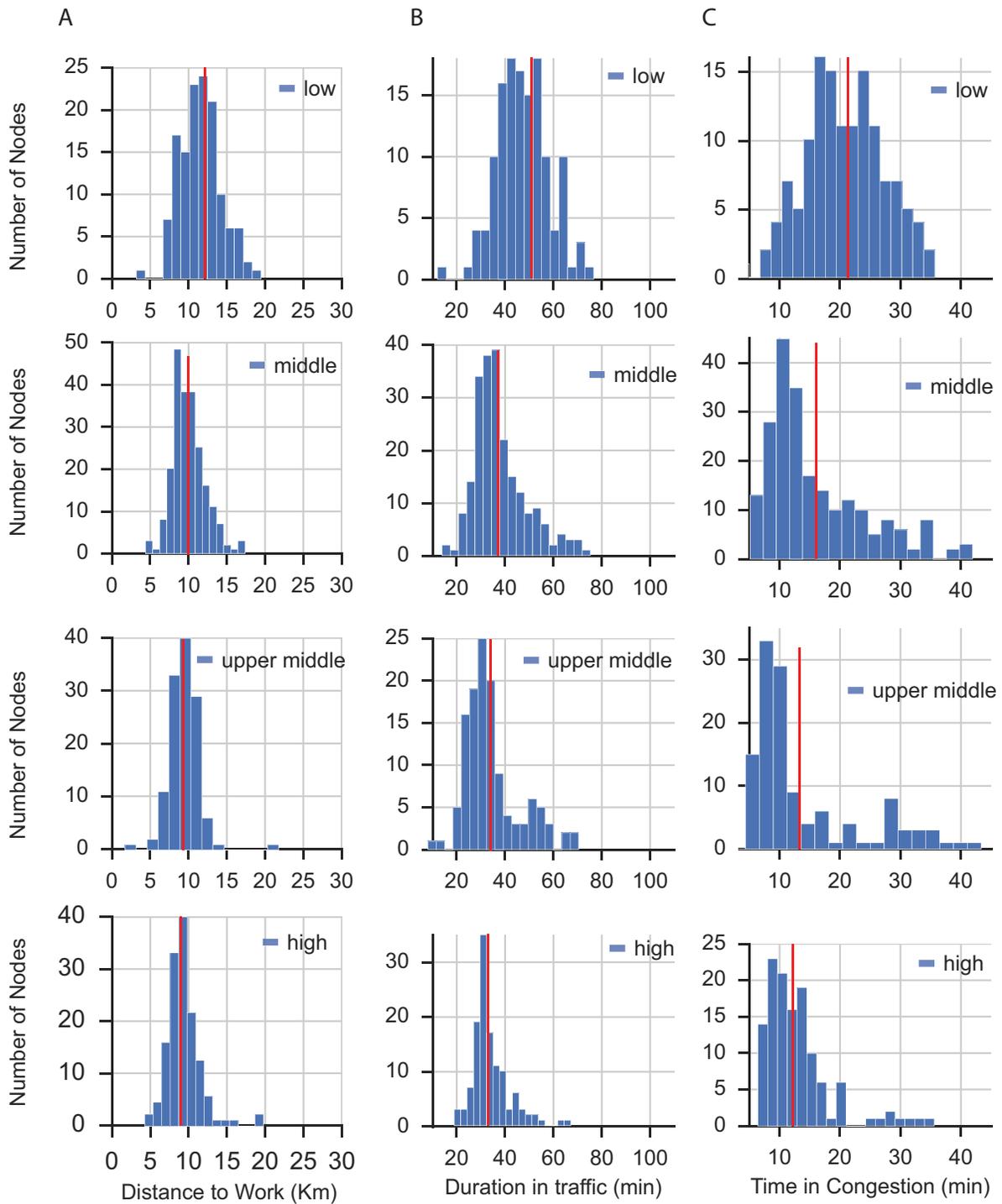


FIGURE 8 Travel by income group for commuters in a sampled network containing 20,000 travelers from each income bracket. A. Distance B. Duration of trips in congested travel time (calculated with the Google Maps api) C. Time spent in congestion.

1 particular, the clear differences in the time spent in congestion by income group, suggest that
2 travel information based on mobile phones can be used to measure the effects of transportation
3 interventions to alleviate congestion for the most vulnerable sectors. Interesting avenues for future
4 research related to this work, include the study of the characteristics of the social network in
5 relation to the economic well-being at the urban scale. Also it is interesting to measure metrics of
6 spatial exposure among different groups and how these affect the economic complexity of various
7 cities. In the ideal scenario, we can learn types of urban organizations that promote better outcomes
8 for their inhabitants. In the new data-rich reality of cities, deeper insights into their social and
9 spatial connections will help make the places we live more sustainable, efficient and equitable.

10 ACKNOWLEDGMENTS

11 We thank Bradley Sturt, Nuria Oliver, Alvaro Ramirez Suarez, Gonzalo Durban Diez and Ricardo
12 Hausmann for enlightening discussions during the design and motivation of this work. Data on
13 the model based on the travel diary of Bogota was kindly provided by Laura Lotero. The research
14 reported herein was funded in part by the Interamerican Development Bank, the MIT-Portugal pro-
15 gram, the Samuel Tak Lee Real Estate Entrepreneurship Laboratory at MIT, DOT via the program
16 New England UTC 25 and the Center for Complex Engineering Systems (CCES) at KACST.

17 REFERENCES

- 18 [1] Lora, E. et al., *The quality of life in Latin American cities: Markets and perception*. World
19 Bank Publications, 2010.
- 20 [2] Uribe-Mallarino, C., Estratificación social en Bogotá: de la política pública a la dinámica de
21 la segregación social. *universitas humanística*, Vol. 65, No. 65, 2008.
- 22 [3] Alzate, M. C., *La estratificación socioeconómica para el cobro de los servicios públicos*
23 *domiciliarios en Colombia: ¿Solidaridad o focalización?* CEPAL, 2006.
- 24 [4] Benabou, R., Inequality and growth. In *NBER Macroeconomics Annual 1996, Volume 11*,
25 MIT Press, 1996, pp. 11–92.
- 26 [5] Stiglitz, J. E., *The price of inequality: How today's divided society endangers our future*,
27 2012.
- 28 [6] Roback, J., Wages, rents, and the quality of life. *The Journal of Political Economy*, 1982, pp.
29 1257–1278.
- 30 [7] Blomquist, G. C., Measuring quality of life. *A companion to urban economics*, 2006, pp.
31 483–501.
- 32 [8] Chen, C., J. Ma, Y. Susilo, Y. Liu, and M. Wang, The promises of big data and small data for
33 travel behavior (aka human mobility) analysis. *Transportation Research Part C*, 2016.
- 34 [9] Çolak, S., A. Lima, and M. C. González, Understanding congested travel in urban areas.
35 *Nature Communications*, Vol. 7, 2016.
- 36 [10] Çolak, S., L. P. Alexander, B. G. Alvim, S. R. Mehndiretta, and M. C. González, Analyzing
37 cell phone location data for urban travel: current methods, limitations and opportunities. In
38 *Transportation Research Board 94th Annual Meeting*, 2015, 15-5279.

- 1 [11] Green, R. K., K. D. Vandell, et al., *Increasing Homeownership Opportunities through Modi-*
2 *fication of the IRS Rules Affecting Owner-Occupied Housing*. University of Wisconsin Center
3 for Urban Land Economic Research, 1995.
- 4 [12] Lotero, L., A. Cardillo, R. Hurtado, and J. Gómez-Gardeñes, Several multiplexes in the same
5 city: The role of socioeconomic differences in urban mobility. In *Interconnected Networks*,
6 Springer, 2016, pp. 149–164.
- 7 [13] Toole, J. L., S. Colak, B. Sturt, L. P. Alexander, A. Evsukoff, and M. C. González, The path
8 most traveled: Travel demand estimation using big data resources. *Transportation Research*
9 *Part C: Emerging Technologies*, 2015.
- 10 [14] Stopher, P. R. and S. P. Greaves, Household travel surveys: Where are we going? *Trans-*
11 *portation Research Part A: Policy and Practice*, Vol. 41, No. 5, 2007, pp. 367 – 381, bridging
12 Research and Practice: A Synthesis of Best Practices in Travel Demand Modeling.
- 13 [15] Blondel, V. D., A. Decuyper, and G. Krings, A survey of results on mobile phone datasets
14 analysis. *EPJ Data Science*, Vol. 4, No. 1, 2015, p. 1.
- 15 [16] Schneider, C. M., V. Belik, T. Couronné, Z. Smoreda, and M. C. González, Unravelling daily
16 human mobility motifs. *Journal of The Royal Society Interface*, Vol. 10, No. 84, 2013, p.
17 20130246.
- 18 [17] Wesolowski, A., T. Qureshi, M. F. Boni, P. R. Sundsøy, M. A. Johansson, S. B. Rasheed,
19 K. Engø-Monsen, and C. O. Buckee, Impact of human mobility on the emergence of dengue
20 epidemics in Pakistan. *Proceedings of the National Academy of Sciences*, Vol. 112, No. 38,
21 2015, pp. 11887–11892.
- 22 [18] Wesolowski, A., N. Eagle, A. J. Tatem, D. L. Smith, A. M. Noor, R. W. Snow, and C. O.
23 Buckee, Quantifying the impact of human mobility on malaria. *Science*, Vol. 338, No. 6104,
24 2012, pp. 267–270.
- 25 [19] Tizzoni, M., P. Bajardi, A. Decuyper, G. K. K. King, C. M. Schneider, V. Blondel,
26 Z. Smoreda, M. C. González, and V. Colizza, On the use of human mobility proxies for
27 modeling epidemics. *PLoS Computational Biology*, Vol. 10, No. 7, 2014, p. e1003716.
- 28 [20] Lu, X., L. Bengtsson, and P. Holme, Predictability of population displacement after the 2010
29 Haiti earthquake. *Proceedings of the National Academy of Sciences*, Vol. 109, No. 29, 2012,
30 pp. 11576–11581.
- 31 [21] Alexander, L., S. Jiang, M. Murga, and M. C. González, Origin–destination trips by purpose
32 and time of day inferred from mobile phone data. *Transportation Research Part C: Emerging*
33 *Technologies*, 2015.
- 34 [22] Iqbal, M. S., C. F. Choudhury, P. Wang, and M. C. González, Development of origin–
35 destination matrices using mobile phone call data. *Transportation Research Part C: Emerging*
36 *Technologies*, Vol. 40, 2014, pp. 63–74.

- 1 [23] Wang, P., T. Hunter, A. M. Bayen, K. Schechtner, and M. C. González, Understanding road
2 usage patterns in urban areas. *Scientific reports*, Vol. 2, 2012.
- 3 [24] of Statistics (Colombia), N. A. D., *METODOLOGÍA DE ESTRATIFICACIÓN URBANA*
4 *TIPO 1*, 2004.
- 5 [25] Jeffrey C. Moore, L. L. S. and J. Edward J. Welniak, Income Measurement Error in Surveys:
6 A Review. *Journal of Official Statistics*, Vol. 16, No. 4, 2000, pp. 331–361.
- 7 [26] Secretaria Distrital de Movilidad, ., Informe de indicadores Encuesta de Movilidad de Bogotá,
8 2011.
- 9 [27] Hariharan, R. and K. Toyama, Project Lachesis: Parsing and Modeling Location Histories.
10 In *Geographic Information Science* (M. Egenhofer, C. Freksa, and H. Miller, eds.), Springer
11 Berlin Heidelberg, Vol. 3234 of *Lecture Notes in Computer Science*, 2004, pp. 106–124.
- 12 [28] Saberi, M., H. S. Mahmassani, D. Brockmann, and A. Hosseini, A complex network perspec-
13 tive for characterizing urban travel demand patterns: graph theoretical analysis of large-scale
14 origin–destination demand networks. *Transportation*, 2016, pp. 1–20.
- 15 [29] Eagle, N., M. Macy, and R. Claxton, Network Diversity and Economic Development. *Sci-*
16 *ence*, Vol. 328, No. 5981, 2010, pp. 1029–1031.
- 17 [30] Pappalardo, L., M. Vanhoof, L. Gabrielli, Z. Smoreda, D. Pedreschi, and F. Giannotti, An
18 analytical framework to nowcast well-being using mobile phone data. *International Journal*
19 *of Data Science and Analytics*, 2016, pp. 1–18.