

# Heterogeneous Preferences for Neighborhood Amenities: Evidence from GPS Data\*

Cody Cook<sup>†</sup>

January 2023

## Abstract

I study how preferences for neighborhood amenities vary by income. Using data on over 150 million visits to restaurants, shops, personal services, and entertainment places, I estimate a model of demand for amenities. I find that higher and lower income urban residents have heterogeneous preferences for individual establishments, which often vary systematically along observable dimensions such as category, brand, and price level. Using the location and estimated quality of each establishment, I construct an aggregate Neighborhood Amenity Quality Index (NAQI) that measures the value of each neighborhood's overall access to amenities. Despite the heterogeneity in establishment-level preferences, neighborhood-level preferences exhibit strong positive correlation; higher and lower income residents generally agree on which neighborhoods have the best overall access to amenities. Densely populated neighborhoods close to the urban core benefit from residential agglomeration forces and have especially high quality access to amenities. Conditional on population density, neighborhoods with better amenity access tend to be richer, more educated, and have more expensive rents.

---

\*I am grateful for valuable feedback from Lanier Benkard, Nick Bloom, Victor Couture, José Ignacio Cuesta, Lindsey Currier, Rebecca Diamond, Liran Einav, Matthew Gentzkow, Ed Glaeser, Jessie Handbury, Pearl Li, Neale Maloney, Peter Reiss, Paulo Somaini, Winnie van Dijk, Shoshana Vasserman, and many seminar participants. This project began while I was a contractor for the urban data platform company Replica, which provided the data and necessary computational resources; Replica did not ask for nor receive any editorial oversight in either the topic selection or the analyses conducted. Steven Kim, Dave Lawlor, Kiran Jain, Aude Marzuoli, and Alexei Pozdnoukhov at Replica all provided valuable assistance.

<sup>†</sup>Stanford GSB. Email: codycook@stanford.edu.

# 1 Introduction

Cities have become increasingly segregated by income over the past 30 years (Reardon et al., 2018), leading to rising concerns about the distributional consequences of income segregation. Recent research has shown how neighborhood amenities can amplify within-city spatial sorting and inequality by endogenously responding to changes in neighborhood demographics, which can encourage additional in-migration and reduce a neighborhood’s appeal for incumbent residents (Couture et al., 2019; Almagro and Dominguez-Iino, 2022). The role of certain categories of amenities, such as restaurants and shops, has been prominently featured in the discourse of neighborhood gentrification.<sup>1</sup> In response, policymakers often target the composition of a neighborhood’s amenities through policies such as zoning restrictions, the designation of landmarks and cultural heritage districts, and bans on chain stores (Klopack, 2021).

The extent to which neighborhood amenities can amplify spatial sorting and inequality will depend on the correlation of preferences for amenities (Waldfogel, 2010; Almagro and Dominguez-Iino, 2022). When preferences for amenities are more aligned across the income distribution, new amenities that enter a neighborhood are likely to have broad appeal. When preferences are less aligned, however, amenities that enter a gentrifying neighborhood to target new, wealthier residents may crowd out those amenities enjoyed by the incumbent, lower income residents. A challenge for policymakers interested in ensuring that neighborhood amenities provide value to all residents is identifying the types of amenities for which residents have more aligned preferences.

In this project, I study how the preferences for restaurants, shops, personal services, and entertainment places vary by income. To separate preferences for amenities from access to these amenities, I build and estimate a model of demand for individual establishments within a city. Agents in the model receive stochastically arriving opportunities to visit some amenity and then choose to frequent a specific establishment (or their outside option) based on the location and quality of each establishment. Decisions are made hierarchically: first agents choose a category (e.g., shops), then a subcategory (e.g., grocery stores), and finally a specific establishment.<sup>2</sup>

I estimate the model using large-scale data on the locations of over 1.2 million amenities in the 100 largest US cities, combined with data on individual visits to these locations by over 15 million GPS-enabled devices during 2019. Amenities are partitioned into the four main categories—restaurants, shops, personal services, and entertainment—and 45 subcategories. For each device in the data, I use parcel-level data on their place of residence to estimate whether they are above the city’s median household income (‘higher income’) or below median income (‘lower income’). I do not observe any prices (and some amenities, such as libraries, are free), so I specify the model in terms of driving time from home and work and convert all measures of utility into units

---

<sup>1</sup>See, for example, recent articles in the New York Times (Gordinier, 2016; Kolomatsky, 2020), the Wall Street Journal (Raleigh, 2017; Ukueberuwa, 2020), and the Atlantic (Smith, 2016)

<sup>2</sup>I use ‘establishment’ to refer to the location of any neighborhood amenity, including both private establishments and public amenities such as parks.

equivalent to saving a minute of driving time. The estimated model provides two key ingredients for understanding the value of neighborhood amenities: 1) establishment-level preference estimates for each income group and 2) a framework for valuing a neighborhood's overall access to amenities.

I find that higher and lower income residents' preferences for individual establishments are positively correlated in all 45 subcategories of amenities, but that the degree of correlation varies substantially. For large retail subcategories, such as malls and general merchandise stores, higher and lower income preferences have a correlation coefficient of over 0.92. Within these subcategories, the entry of any establishment is likely to provide similar value to a broad swathe of the income distribution. Meanwhile, preferences for subcategories such as restaurants, barber shops, and gyms are less correlated. Within these subcategories, there is more scope for establishments that enter a neighborhood to tailor to a specific sub-population to the detriment of other residents. For a subset of establishments, I also observe characteristics such as brand affiliation and Yelp price levels. These underlying characteristics and are often predictive of which income group has stronger preferences for an establishment – for example, higher income residents tend to prefer restaurants with higher Yelp prices and cuisines such as New American, while lower income residents often prefer cheaper restaurants and those that are affiliated with a large chain.

I next turn to how preferences for a neighborhood's overall access to amenities vary by income. I define a 'Neighborhood Amenity Quality Index' (NAQI) that captures the value of amenity access for higher and lower income residents of different Census block groups within a city.<sup>3</sup> Neighborhoods with more high quality establishments nearby will have a higher NAQI value. I find that preferences for neighborhoods' overall level of amenity access are far more correlated than preferences for individual establishments – the average within-city correlation in NAQI values is 0.90. Rather than neighborhoods tailored to rich residents and others tailored to poor residents, cities generally contain neighborhoods dense in amenities with sufficient variety to have broad appeal and other neighborhoods with more limited access. The primary determinant of the value of neighborhoods' access to amenities is population density. Dense neighborhoods close to the urban core benefit from the residential agglomeration effects and tend to have a rich set of amenities to explore. Conditional on density, neighborhoods with higher household income, rent, and education tend to have more valuable access to amenities. Despite the high level of correlation in neighborhood-level preferences, some systematic differences do emerge; for example, the relationship between NAQI values and neighborhood characteristics such as rent and education is about twice as strong for higher income residents than for lower income residents.

Overall, the high levels of positive correlation in preferences for neighborhoods suggest that direct effects of tailoring these types of amenities to local demographics will generally be second-order to the total density of nearby amenities. While the short panel of available data preclude studying the joint process of demographic change and amenity entry/exit directly, I use the cross-

---

<sup>3</sup>The estimated NAQI values for each neighborhood are available by request and will be publicly posted in the future.

sectional estimates to simulate a fairly extreme version of gentrification by replacing all of the top 25% establishments for lower income residents with replicas of the top 25% of establishments for higher income residents. In the average city, this has a modest effect on total utility, equivalent to +18.1 minutes of weekly driving time saved for higher income residents and -21.9 minutes for lower income residents (respectively +\$313 and -\$371 annually if residents value their time at \$20/hour). Nonetheless, concurrent changes in rents, other amenities such as school quality, job access, social cohesion, and “neighborhood character”, or substantial re-sorting by residents in response to changes in amenities may still lead to large overall welfare effects.

Earlier work on estimating the value of neighborhood amenities has generally taken one of two approaches. The first approach treats the value of amenities as a residual in a residential choice model, building on the canonical Rosen-Roback model (Rosen, 1979; Roback, 1982).<sup>4</sup> This approach has become especially common in quantitative spatial economics models, reviewed in Redding and Rossi-Hansberg (2017). In one notable example, Ahlfeldt et al. (2015) find that the elasticity of amenity quality to population density is twice that of the elasticity of workplace productivity to density, consistent with the strong relationship between amenity quality and population density that I find in this paper. The second common approach, which I use here, is a ‘bottom-up’ approach that uses data on individual choices to estimate revealed preferences for neighborhood amenities.<sup>5</sup> As two examples, Couture (2016) combines data on the locations of restaurants and the travel times of local residents going out to eat to show that the value of areas dense in amenities is primarily due to variety, rather than shorter travel times, and Davis et al. (2019) uses geolocation data from Yelp reviews to show how both spatial and social frictions lead to ethnic and racial segregation in consumption of restaurants. The bottom-up approach has traditionally been limited by lack of comprehensive data on the locations of amenities matched to individual visits to each location. My use of large-scale GPS data enables measuring the value of neighborhood amenities for many different categories and across all large US cities.

A related line of research has studied how preferences for neighborhood amenities can affect residential sorting. Recent work by Couture et al. (2019), Baum-Snow and Hartley (2020) and Couture and Handbury (2020) finds that the return of higher income, more educated residents to urban cores is due to a rising tendency for these residents to sort towards areas rich in non-tradeable, consumption amenities.<sup>6</sup> These sorting-on-amenities results have two possible underlying sources: differences in how groups of residents value the specific amenities located in urban cores or differences in the overall weight they place on neighborhood amenities when making residential

---

<sup>4</sup>Albouy (2008) find that the Rosen-Roback quality of life measures are positively correlated with observable amenities.

<sup>5</sup>Similar work by Handbury and Weinstein (2015) and Handbury (2021) use data on product-level choices to estimate local price indices.

<sup>6</sup>This more recent work builds on a long history of research into urban amenities, neighborhood change, and gentrification. See, for example, Brueckner et al. (1999); Glaeser et al. (2001); Vigdor et al. (2002); Ellen and O'Regan (2010); McKinnish et al. (2010); Guerrieri et al. (2013); Autor et al. (2017); Lee and Lin (2018); Glaeser et al. (2018); Su (2022); Fogli and Guerrieri (2019); Brummet and Reed (2019)

choices. The results here suggest the latter explanation is more likely as I find that higher and lower income residents *both* value the amenity access of the urban core far more than that of other neighborhoods. The recent trends in sorting, therefore, likely stem from differences in overall willingness to pay for amenity access – while lower income residents still prefer the amenities of the urban core, they are displaced by higher income residents with greater willingness to pay for better access to amenities.

Other research has emphasized how heterogeneity in preferences for amenities combined with endogenous amenities can amplify spatial inequality. As demographics shift, neighborhood amenities will tailor to the tastes of newcomers, which can magnify demographic sorting and reduce the welfare of incumbent residents. This research extends work on spillovers between demographic groups through supply-side responses to local preferences (Waldfogel, 2003, 2010). At the city-level, Diamond (2016) shows that the endogenous response of amenities in high-skill cities fuels increased across-city sorting by skill. More recent papers extend the intuition from Diamond (2016) to show how endogenous amenities can reinforce within-city sorting (Couture et al., 2019; Hoelzlein, 2019; Almagro and Dominguez-Iino, 2022). My estimates of subcategory-specific preference correlations highlight that certain types of amenities with less aligned preferences (e.g., personal services) have more potential to act as amplifiers of spatial inequality than other types with more aligned preferences (e.g., general merchandise).

Finally, this project draws from a developing literature using data from GPS-enabled devices to study questions related to urban mobility. Applications of GPS data include mobility during the COVID-19 pandemic (Chang et al., 2020; Allcott et al., 2020), waiting times at voting polls (Chen et al., 2022), the effect of new transit options on commuting behavior (Gupta et al., 2022), knowledge spillovers between employees of different firms (Atkin et al., 2020), and alternative measures of segregation based on where people spend their time (Caetano and Maheshri, 2019; Athey et al., 2021; Cook et al., 2022). Most related to my own work, Athey et al. (2018) estimate demand for lunch restaurants and Miyauchi et al. (2021) study how the consumption externalities of commute patterns shape the spatial structure of cities. A common downside of GPS data is the lack of individual-level demographics; in this project, I use a novel method for inferring individual-level income by matching devices to their precise home parcel.

## 2 Data

**Locations of neighborhood amenities.** I identify neighborhood amenities using SafeGraph’s Places dataset, which includes the name, location, category, associated brands, and the polygon describing the location’s footprint for a large number of Points of Interest (POIs).<sup>7</sup> I subset to the

---

<sup>7</sup>The data cover only those amenities which have established physical footprints, which excludes many features of a neighborhood that residents may appreciate, such as trees, clean streets, access to quality schools, and public transportation, but encompasses a wide range of both commercial and public locations such as restaurants, shops,

100 most populous Core-Based Statistical Areas (CBSAs) to focus on urban areas. I categorize amenities into four main categories: restaurants, shops, personal services, and entertainment places. Within each category, I assign more granular subcategories based on each establishment’s North American Industry Classification System (NAICS) code. For example, restaurants are further divided into full service, limited service, cafes, and drinking places. Shops are divided into grocery stores, malls, department stores, book stores, and more. Personal services include places such as barber/salons, banks, religious organizations, hospitals, mechanics, and entertainment places include cinemas, parks, gyms, golf courses, performing arts venues, and others. Appendix Table [A1](#) documents the full set of 45 subcategories, the NAICS codes used to identify them, and the number of POIs within each.

**Amenity characteristics from Yelp.** I augment the SafeGraph data using data scraped from Yelp on the average rating, number of reviews, and price level for each establishment that has a Yelp page. The availability of these features varies by amenity category; most restaurants have pages, but other amenities such as parks may not.

**Visits to amenities.** I build a dataset of visits to individual locations during 2019 using a sample of location histories from GPS-enabled devices provided by Replica. The raw data consists of GPS ‘stays,’ each of which includes a unique device ID, entry time, exit time, and the coordinates visited. Home and work locations are assigned for each device-quarter using heuristics on when people are at home and work. At a high-level, a device’s home is defined as the most frequent overnight location and its work is the most frequent non-home daytime location. Stays outside of home and work are matched to SafeGraph POIs using the POI polygons. For POIs contained within other POIs—e.g., a store within a large mall—the ‘parent’ POI is considered the visited location. Appendix Section [A.1](#) contains additional details on the GPS data construction.

**Individual income.** A downside of GPS data is that they do not include any characteristics of the individual. Past research using GPS data has often inferred demographics based on the median household in its home location (e.g., Census block group). However, in mixed-income block groups, using the median income to classify devices as higher or lower income will introduce substantial measurement error. Instead, I combine information on the characteristics of a device’s home parcel, an estimate of the market value of their home, and the income distribution of their home block group to predict whether each device is part of a household that is above the median household income in the CBSA (‘higher income’) or below median income (‘lower income’). I describe the full

---

barbers, gyms, parks, churches, zoos, libraries, dentists, banks, and more. [Abbiasov and Sedov \(2022\)](#) compares the number of establishments in different categories in both Safegraph and in the Census’s County Business Patterns and finds that coverage is similar.

procedure in Appendix Section A.2. In Section 5.4, I show that my results are robust to alternative measures of income.

**Driving times.** I use driving times—rather than crow-flies distance—to measure how accessible a POI is from a resident’s home and work locations. I compute driving times using Graphhopper, a router built on OpenStreetMaps. To reduce the number of routing requests, I pre-compute driving times between all pairs of block groups within each CBSA, using the block group centroids as the start and end points. For the largest 100 CBSAs, this involves computing nearly 750 million total routes.

**Estimation sample.** The final sample consists of 15 million devices living in the largest 100 CBSAs nationwide. Table 1 provides a number of summary statistics for the aggregate data as well as split out for a subset of large CBSAs. For the average device-quarter, I observe 93.2 stays, of which 47.4 are at home, 11.8 are at work, and 15.7 are at some SafeGraph POI. Appendix Table A1 documents the number of visits and quantiles of the driving time distribution for each subcategory of amenity.

**Sample weights.** The GPS data is not a random sample of the US population – it tends to over-sample from slightly lower income, less educated, and less white block groups within a city. I discuss a number of measures of sample quality in Appendix Section A.3. For all analyses, I use sample weights to adjust for non-uniform sampling. I weight each device-quarter by  $\omega_{iq} = N_{g(iq)}^{Total} / N_{g(iq)}^{GPS}$  where  $N_{g(iq)}^{GPS}$  is the number of GPS devices observed in the device’s home block group  $g(iq)$  and  $N_{g(iq)}^{Total}$  is the 2019 5-year ACS population of the block group.

### 3 Consumption and access to local amenities

In this section, I show how higher and lower income residents differ in both their access to and levels of consumption of different categories of neighborhood amenities.

First, higher income residents live near substantially fewer amenities of all categories. To show this, I regress the log number of establishments within driving radii of 5, 15, and 30 minutes from each individual’s home on whether an individual is higher income. The average higher income resident, controlling for their home CBSA, has 22% fewer shops and restaurants, 15% fewer personal services, and 6% fewer entertainment places within a 15 minute drive from home (Figure 1 Panel a). While many higher income residents have begun to move towards city centers in recent decades, the majority continue to live in less dense suburbs with worse access to amenities (Couture et al., 2019). Appendix Figure C.2 repeats the regressions, but controlling for the log population density of device’s home block group. Controlling for density, higher income residents have 5-20% *more*

Table 1: GPS device sample

		Example CBSAs				
	All	Chicago	Bay Area, CA	Los Angeles	Houston	Phoenix
Num. devices	15,090,847	634,050	316,709	863,717	525,229	356,265
Num. device-quarters	24,698,171	1,048,719	517,750	1,377,706	852,045	569,294
% higher income	52.71	52.92	53.78	53.59	50.93	51.69
% with work	67.29	67.65	66.43	63.67	68.49	62.37
Avg. quarterly stays	93.21	94.53	91.98	90.6	98.58	87.32
	[85.13]	[85.52]	[83.59]	[85.33]	[88.88]	[81.36]
Avg. quarterly stays at POI	15.66	16.35	17.02	15.52	16.92	15.71
	[23.51]	[23.8]	[24.15]	[23.63]	[24.74]	[24.69]
Avg. quarterly stays at home	47.37	48.36	47.23	47.08	48.18	46.55
	[38.61]	[38.95]	[38.36]	[39.2]	[38.67]	[38.62]
Avg. quarterly stays at work	11.82	12.03	11.87	11.31	12.66	10.45
	[17.49]	[17.48]	[17.72]	[17.71]	[18.41]	[16.79]

*Note:* This table documents the size of the sample and the number of stays observed for the average device-quarter. Stays at POI are those stays that are matched to an exact location in the SafeGraph data. Stays at home and work are the average number of stays labeled as at the device’s home or work location. Standard deviations are presented in brackets. All statistics are weighted by the device’s sample weight.

amenities within a 15 minute drive, although continue to have fewer amenities within a 5 minute drive.

Second, higher income residents tend to work near more amenities, although the gap is much smaller than with amenity access from home. The average higher income resident has access to 2-9% more establishments of each category within a 15 minute drive of their workplace (Figure 1 Panel b). This is again closely tied to density; while higher income residents tend to live in the suburbs of cities, they commute to denser areas to work.

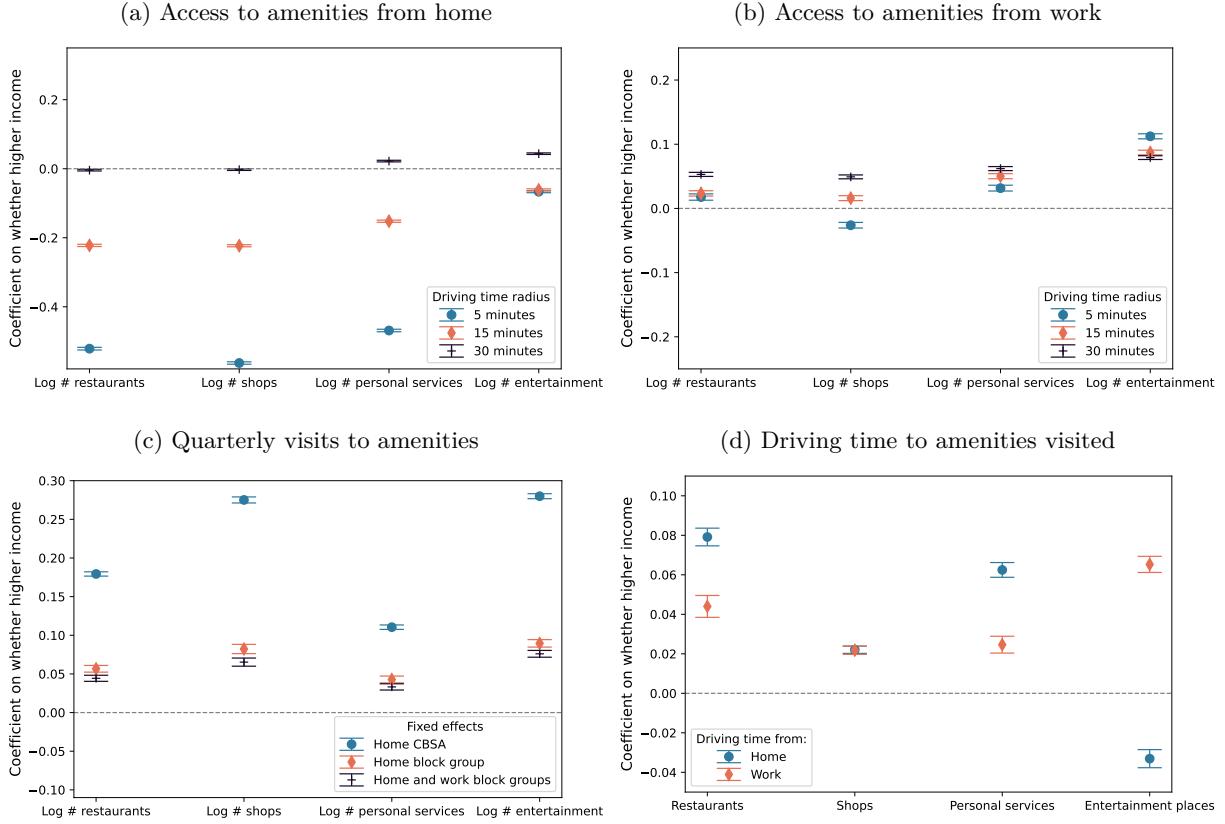
Third, despite worse access to amenities from home and only marginally better access from work, higher income residents visit far more amenities each quarter (Figure 1 Panel c). On average, a higher income resident will visit 14% more restaurants, 23% more shops, 9% more personal services, and 22% more entertainment places than the average lower income resident. Even when controlling for both home and work block groups, higher income residents continue to visit 2-6% more amenities of each category each quarter, suggesting both access and differences in preferences play a role.<sup>8</sup>

Finally, in order to visit different amenities, higher income residents will travel an average of 2-8% further from home and work than lower income residents each time (Figure 1 Panel d). Even when controlling for home block group, higher income residents continue visit amenities that are 1-2% further from home and work than lower income amenities (Appendix Figure C.2).

<sup>8</sup>This gap in number of visits even when controlling for home and work locations also helps motivate my use of an individual-level measure of income, rather than using the median household income of their home block group.

Appendix Figure C.3 repeats these analyses for each of the 45 subcategories of amenities. The largest difference between higher income and lower income consumption are for full-service restaurants, malls, parks, golf courses, and fitness centers. The only categories that higher income residents tend to have *better* access to from home are golf courses, drycleaning/laundry, dentists, and other health practitioners (e.g., chiropractors). The broad trends are largely similar across both higher level categories as well as subcategories of amenities – despite having worse access to these amenities from home and only slightly better access from work, higher income residents consume far more amenities each quarter.

Figure 1: Neighborhood amenities: access and consumption



*Note:* This figure plots coefficients from a series of regression where the right hand side is an indicator for whether a device is above median income. Each regression includes fixed effects for the CBSA. Each marker in a plot covers a 95% confidence interval, with standard errors clustered at the device level. Panel (a) and (b) are at the device-quarter level and regress the log number of establishments within different lengths of driving time. The outcomes for Panel (c) is the log number of visits to establishments of each category that quarter. Panel (d) is at the visit level and regresses the log of driving times from home and work on whether the device is higher income. To handle zeros, I use an inverse hyperbolic sine transformation instead of log, although I label it as ‘log’ for brevity.

## 4 Modeling amenity choice

The descriptive facts highlight that higher and lower income residents of a city consume different levels and categories of neighborhood amenities. However, these descriptives conflate preferences and access; while consumption varies by income, so does the level of access. To separate preferences from access, I estimate a choice model of amenity consumption. The model provides a unified framework for studying both correlations in preferences for individual establishments as well as the value of a neighborhood's overall access to amenities.

More concretely, I model demand for amenities as a three-level nested logit, illustrated in Figure 2. Opportunities to consume some amenity arrive stochastically according to a Poisson arrival process with an arrival rate that can vary by time of day and income. This arrival process mimics different levels of external time constraints; a lower arrival rate corresponds to having less availability to visit amenities. Upon receiving a choice opportunity, an individual first chooses between visiting a restaurant, shop, personal service, entertainment place, or their outside option. Second, she chooses a subcategory within the chosen main category. For example, within shops she may choose between malls, grocery stores, book stores, and more. Finally, she chooses a specific establishment within a subcategory. Her choice will depend on the time of day, the driving times from home and work of different establishments, the type and quality of each establishment, and an idiosyncratic component.

When presented with a choice opportunity at some time of week  $t$ , the probability that an individual  $i$  chooses establishment  $j$  can be decomposed into three probabilities, corresponding to each level of the model:

$$P_{itj} = P_{itj|B_m} * P_{itB_m|B_n} * P_{itB_n} \quad (4.1)$$

where  $B_m$  is a subcategory nest and  $B_n$  is a main category nest.  $P_{itj|B_m}$ ,  $P_{itB_m|B_n}$ , and  $P_{itB_n}$  each take the form of a logit. Below, I describe in detail how each level is modeled, starting from the lowest level.

### 4.1 Lowest level: choice of a specific establishment

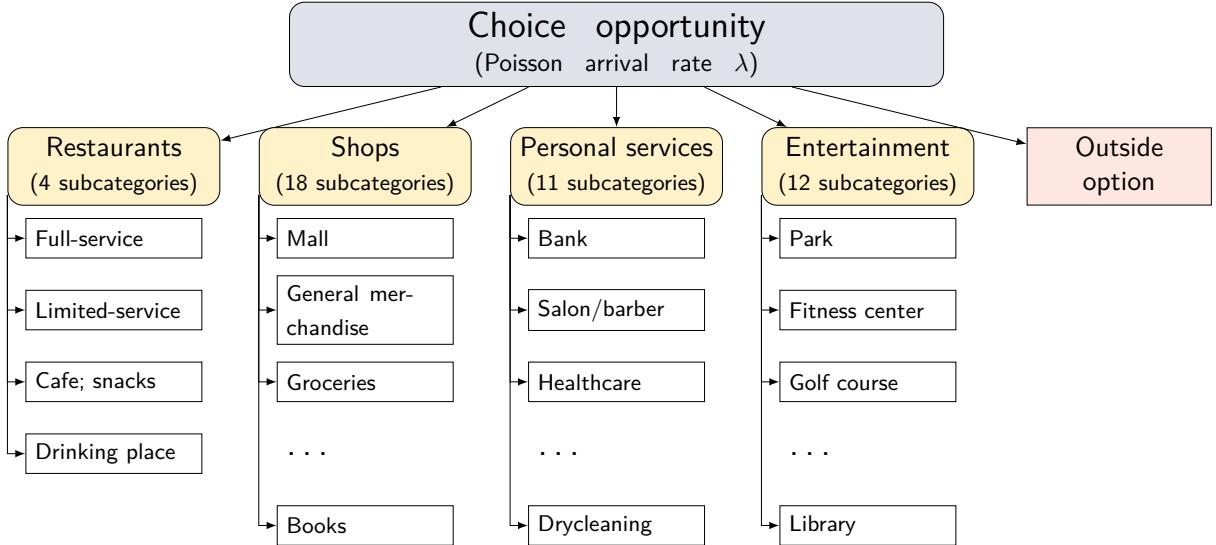
Conditional on choosing to consume within subcategory  $B_m$ , an individual next selects a specific establishment to visit. Her choice set consists of all of the subcategory's establishments within the CBSA. The indirect utility from a given establishment  $j$  at time of week  $t$  is

$$u_{ijt} = \gamma_{k(i)j} - \kappa_{k(i)t}^H d_{ij}^H - \kappa_{k(i)t}^W d_{ij}^W + \epsilon_{ijt} \quad (4.2)$$

$$= \delta_{ijt} + \epsilon_{ijt} \quad (4.3)$$

where the individual's type is denoted  $k(i)$ ,  $\gamma_{k(i)j}$  are establishment-type taste parameters, and  $d_{ij}^H$  and  $d_{ij}^W$  are driving times from home and work with corresponding disutilities  $\kappa_{k(i)t}^H$  and  $\kappa_{k(i)t}^W$  that

Figure 2: Model structure



vary by time of week and individual type. The model is defined separately for each subcategory  $B_m$ , although I omit any subcategory subscripts for clarity. I discretize time of week  $t$  into six bins corresponding to (weekday, weekend)  $\times$  (morning, afternoon, evening).<sup>9</sup> The idiosyncratic component ( $\epsilon_{ijt}$ ) is redrawn each choice opportunity, so there is no persistence in idiosyncratic preferences across choice opportunities. Following McFadden et al. (1973), I assume that these errors are distributed according to a type 1 extreme value distribution, which gives the usual conditional choice probabilities of consuming at a specific establishment within subcategory  $B_m$ :

$$P_{ijt|B_m} = \frac{\exp(\delta_{ijt})}{\sum_{l \in B_m} \exp(\delta_{ilt})} \quad (4.4)$$

Unlike standard demand settings, I do not observe any prices. Instead, the ‘cost’ of consuming at any given establishment is the driving time from home and, where applicable, work. As such, the  $\gamma_{k(i)j}$  capture the ‘price inclusive’ value of establishment  $j$  to an average user of type  $k$ . These values are identified from 1) the frequency at which an establishment is chosen from the set of alternatives and 2) the driving time individuals were willing to travel in order to visit. Intuitively, if an establishment is frequently visited by residents who live and work far away, it must have a high  $\gamma_{k(i)j}$  term to rationalize this behavior.

I normalize the model using the disutility of driving by fixing the disutility of driving an extra minute from home on weekday evenings to be  $-1$  for all subcategories. This has two advantages. First, it allows the scale of the idiosyncratic component of utility to vary across different subcategories of amenities. Second, it converts the utility of visiting a given establishment into interpretable

---

<sup>9</sup>For times of week, morning is 2am to 11am, afternoon is 11am to 5pm, and evening is 5pm to 2am the next day. Friday evenings are considered weekend evenings and Sunday evenings are considered weekday evenings.

units. I will use  $\tilde{\gamma}_{kj}$  to denote the normalized establishment-type taste parameters for individuals of type  $k$  consuming at establishment  $j$  and  $\tilde{\delta}_{ijt}$  to denote the normalized mean utility for individual  $i$  at time of week  $t$ . These are given by

$$\tilde{\gamma}_{kj} \equiv \frac{\gamma_{kj} - \bar{\gamma}_k}{\kappa_{k0}^H}, \quad \tilde{\delta}_{ijt} \equiv \tilde{\gamma}_{k(i)j} - \left( \frac{\kappa_{k(i)t}^H d_{ij}^H + \kappa_{k(i)t}^W d_{ij}^W}{\kappa_{k(i)0}^H} \right)$$

where I also re-center the establishment parameters to be relative the median establishment for this subcategory and individual type ( $\bar{\gamma}_k$ ). Therefore, the difference  $\tilde{\gamma}_{kq} - \tilde{\gamma}_{kr}$  for two establishments  $q$  and  $r$  corresponds to the difference in utility from consuming at establishment  $q$  instead of  $r$ , with units that correspond to the number of minutes an individual of type  $k$  would be willing to drive in order to go to establishment  $q$  instead of  $r$ .

## 4.2 Upper levels: choice of a category and subcategory

In the upper levels of the model, when an individual is presented with a choice opportunity she chooses a category of amenity (or the outside option) and then a subcategory within her chosen category. The relatively probability of choosing to consume within a given subcategory will depend on the quality of her access to establishments of that subcategory based on her home and work locations and the estimated establishment-level parameters for her income-type.

I specify the indirect utility of consuming some alternative in subcategory  $B_m$  as

$$\begin{aligned} U_{itB_m} &= W_{itB_n} + Y_{itB_m} + \epsilon_{itB_m} \\ &= \nu_{k(i)tB_n} + \mu_{k(i)B_m} + \beta_{k(i)} \text{DTE}_{itB_m} + \epsilon_{itB_m} \end{aligned} \tag{4.5}$$

where the error terms follow the generalized extreme value distribution corresponding to a nested logit model (Ben-Akiva, 1973),  $W_{itB_n} = \nu_{k(i)tB_n}$  is the component of utility from consuming at category  $B_n$  and is parameterized as a set of category intercepts that vary by time of week, and  $Y_{itB_m} = \mu_{k(i)B_m} + \beta_{k(i)} \text{DTE}_{itB_m}$  is the component of utility from consuming at subcategory  $B_m$ . I parameterize  $Y_{itB_m}$  to include a set of subcategory intercepts,  $\mu_{k(i)B_m}$ , and an aggregate measure of ‘driving time equivalents’ (DTE) at level of the subcategory  $B_m$  at time of week  $t$ .

The aggregated DTEs provide a link between the establishment-level estimates in the lower level and the upper levels of the model. These are similar to an inclusive value – they capture the expected maximal draw from establishments in subcategory  $B_m$  for individual  $i$  at time  $t$ , normalized to be in units of minutes of driving time. They are computed as

$$\text{DTE}_{itB_m} = \sigma_{k(i)B_m} \log \left( \sum_{j \in B_m} \exp(\tilde{\delta}_{ij}) \right) \tag{4.6}$$

where  $\sigma_{kB_m}$  is the scale parameter of idiosyncratic component of utility for establishments within subcategory  $B_m$ . The coefficient  $\beta_k$  in Equation 4.5 captures the value of a minute of driving time, through the normalization of utility into these units for each subcategory. In general, the DTEs will be higher for individuals who live and/or work in areas with better access to high quality establishments within this subcategory.

Conditional on receiving a choice opportunity, the probability of consuming at some establishment within subcategory  $B_m$  of main category  $B_n$  is given by

$$\begin{aligned} P_{itB_m} &= P_{itB_m|B_n} \cdot P_{itB_n} \\ &= \left( \frac{\exp(Y_{itB_m}/\rho_{k(i)B_n})}{\sum_{B_\ell \in B_n} \exp(Y_{itB_\ell}/\rho_{k(i)B_n})} \right) \left( \frac{\exp(W_{itB_n} + \rho_{k(i)B_n} IV_{itB_n})}{1 + \sum_{B_h \in B} \exp(W_{itB_h} + \rho_{k(i)B_h} IV_{itB_h})} \right) \end{aligned} \quad (4.7)$$

where the utility of the outside option is normalized to 0,  $B$  is the set of categories,  $\rho_{k(i)B_n}$  measures the degree of independence in unobserved utility for subcategories  $B_m$  within category  $B_n$ , and  $IV_{itB_n}$  are the inclusive values that link the two upper levels. The inclusive values are computed as

$$IV_{itB_n} = \log \left( \sum_{B_m \in B_n} \exp(Y_{itB_m}/\rho_{k(i)B_n}) \right)$$

### 4.3 Neighborhood Amenity Quality Index (NAQI)

The upper levels of the model provide a framework for measuring the overall value of a neighborhood's amenity access, based on the driving times to different establishments for residents of the neighborhood and the estimated preferences for these establishments. In this section, I define a 'Neighborhood Amenity Quality Index' (NAQI) that captures the value of each neighborhood's access to amenities relative to the median neighborhood for each income group, with units that correspond to weekly minutes of driving time.

For a device of type  $k$  living in Census block group  $g$ , the expected value of receiving a choice opportunity at time of week  $t$  is given by

$$EV_{gtk} = \log \left( 1 + \sum_{B_n \in B} \left[ \sum_{B_m \in B_n} \exp \left( \frac{W_{tkB_n} + Y_{gtkB_m}}{\rho_{kB_n}} \right) \right]^{\rho_{kB_n}} \right) \quad (4.8)$$

where I overload notation slightly to let  $W_{tkB_n}$  and  $Y_{gtkB_m}$  represent the portion of utility from the category and subcategory (respectively) at time of week  $t$  for a consumer of type  $k$  who lives in block group  $g$  and has no workplace.

To construct the NAQI values for each neighborhood, I aggregate  $EV_{gtk}$  across all times of week using the estimated weekly arrival rates of choice opportunities. The overall level of utility in Equation 4.8 cannot be determined, so I normalize each neighborhood's value relative to the

median neighborhood for each type in each CBSA. To convert the measure into interpretable units, I standardize using the coefficient on DTEs for a given type:

$$\text{NAQI}_{gk} = \frac{1}{\beta_k} \sum_t [\tilde{\lambda}_{kt} (\text{EV}_{gk t} - \text{EV}_{0k t})] \quad (4.9)$$

where  $\text{EV}_{0k t}$  is the expected value of amenity access for the median neighborhood,  $\beta_k$  is the estimated value of a DTE for group  $k$  from Equation 4.5, and  $\tilde{\lambda}_{kt}$  is the expected number of weekly arrivals for each time of week.<sup>10</sup>

#### 4.4 Estimation

The model is computationally difficult to estimate due to the number of observed choices, the size of the choice sets, and the corresponding number of establishment-level parameters. In this section, I provide an overview of the steps taken to make estimation tractable; Appendix B includes additional details.

I estimate the model for each CBSA and each device type separately.<sup>11</sup> For each of CBSA and device type, I estimate the model sequentially, beginning with estimation of the establishment-level parameters for each separate subcategory. Computation time of the lower level is superlinear in the number of choices within the choice set; sequential estimation allows me to partition the choice sets into just the subset of establishments within a subcategory. Each lower level is a standard logit, which I estimate by maximizing the log-likelihood of observing the set of amenity visits. The likelihood contribution of any single visit is given by the log of  $P_{ijt|B_m}$ , defined in Equation 4.4.

Given the establishment-level estimates ( $\hat{\gamma}$ ) and driving time disutilities ( $\hat{\kappa}$ ), I then jointly estimate the upper levels of the model using aggregate data on the number of visits to different subcategories by each device-quarter. For each individual  $i$  in quarter  $q$  at time of week  $t$ , I observe a vector  $\vec{n}_{iqt}$  of total visits to each subcategory  $B_m$ . The panel nature of the GPS data is necessary to identify the arrival rates of choice opportunities ( $\lambda_{kt}$ ) or, equivalently, the share of choice opportunities for which individuals select the outside option. Given the structure imposed by the nested logit, the arrival rates are pinned down by the cross-category elasticities.<sup>12</sup> The likelihood contribution of observing the vector  $\vec{n}_{iqt}$  for arrival rate  $\lambda_{kt}$  is the sum of the probability

---

<sup>10</sup>In estimation, I use a daily arrival rate  $\lambda_{kt}$  – I multiply this by five for weekdays and two for weekends to get the weekly arrival rate.

<sup>11</sup>I include a buffer region of residents and establishments located within 10 miles of the CBSA’s border. Without this buffer region, neighborhoods close to the border would look mechanically worse than those closer to the center, as many of their nearby amenities would be outside of the CBSA.

<sup>12</sup>To build intuition for this, consider two possible worlds with 10 observed visits to restaurants – one in which choice opportunities are plentiful (i.e. large outside option share) and one in which choice opportunities are infrequent (i.e. small outside option share). Across individuals, for an increase in the DTEs of restaurants holding fixed the value of other categories, the amount of cross-substitution from other categories is informative about which world we are in. When the outside share is small, any increase in restaurant consumption is likely to come via a reduction in consumption of shops, personal services, and entertainment. Meanwhile, if the outside share is large, then cross-substitution from these other categories will be more muted.

of a given number of choice opportunities multiplied by the conditional probability of observing  $\vec{n}_{iqt}$  given that number of choices:

$$P[\vec{n}_{iqt}] = \sum_{a=0}^{\infty} P[a \text{ arrivals} | \lambda_{k(i)t}] * P[\vec{n}_{iqt} | a] \quad (4.10)$$

I derive the full expression in Appendix Section B.3, which I then use to maximize the log-likelihood of observing the full set of choice vectors,  $\{\vec{n}_{itq}\}$ , for all individuals, quarters, and times of week.

Finally, to implement these estimation steps at-scale I take advantage of a relatively new computational framework from the machine learning world called PyTorch, which is especially well-suited to settings where both the data and parameter space are large (Paszke et al., 2017). I provide an overview the advantages of using PyTorch in Appendix Section B.1. The combination of PyTorch, partitioning estimation into individual CBSAs and income groups, and estimating the model sequentially makes estimation feasible for large cities with >50,000 establishments. Using a single GPU, it takes about 4-6 hours to estimate parameters for all device types in a large CBSA such as Chicago. For each city and income group, the estimation procedure provides estimates for each establishment ( $\vec{\gamma}$ ), driving time disutilities ( $\vec{\kappa}$ ), the utility of DTEs for each subcategory and time of week ( $\beta$ ), scale parameters ( $\vec{\rho}$ ), category intercepts ( $\vec{\nu}$ ), subcategory intercepts ( $\vec{\mu}$ ), and arrival rates ( $\vec{\lambda}$ ).

## 4.5 Discussion

The model makes a number of simplifications in the pursuit of tractability at scale. The first major simplification is to limit preference heterogeneity to just two types—higher and lower income residents—rather than more flexible specifications such as within-type random coefficients.

The second simplification is to use only the driving times from home and work and consider each trip in isolation, rather than as a part of a larger chain of trips. A more realistic approach might be to model an individual’s choice of a full sequence of stays during a day, allowing the driving time to be amortized over multiple stops at POIs within a single chain of trips. However, devices in the GPS data are generally not observed for the entirety of the day and chains of trips to more than one POI are rarely observed (see Appendix Section A.1.4 for more details on trip chains). In contrast, Miyauchi et al. (2021) use a small sample of GPS devices in Japan with full coverage to estimate an activity-based model of trip chains.<sup>13</sup> In a robustness check in Section 5.4, I repeat the estimation using only trips that start and end at home and so are not part of a longer trip chain.

Finally, while I parameterize the model in terms of driving time and use a minute of driving as a numeraire, not all trips to amenities are by car and car usage may vary by income, city,

---

<sup>13</sup>Relihan (2022) also explicitly models trip chains to study the effect of online retail on locations such as coffee shops, using data on credit card purchases for estimation.

and destination. As such, the minutes of driving time are better interpreted as an approximation for the time it would take to get between home/work and an amenity, similar to—but hopefully more representative—than the commonly used crow-flies distance. Appendix Figure C.1 plots the normalized disutilities to driving time by time of week. On weekdays, both higher and lower income devices become more sensitive to driving further from home and less sensitive to driving further from work as the day progresses. On weekends, individuals effectively only care about driving time from home.

The primary benefit of these simplifications is to allow estimation of the full model of amenity consumption for all establishments and neighborhoods in the largest 100 CBSAs. In total, I estimate how higher and lower income residents value over 1 million establishments and over 100,000 neighborhoods.

## 5 Results

### 5.1 Who likes which amenities?

#### 5.1.1 Relationship between establishment preferences and observables

In this section, I discuss the relationship between an establishment’s driving time equivalents (DTEs) for each device type and its observable characteristics. While the model does not rely on any characteristics of establishments for estimation, for a limited set of establishments I can observe potentially informative characteristics such as a restaurant’s cuisine or a store’s brand.

To study the relationship between preferences and establishment observables, I run regressions of the form

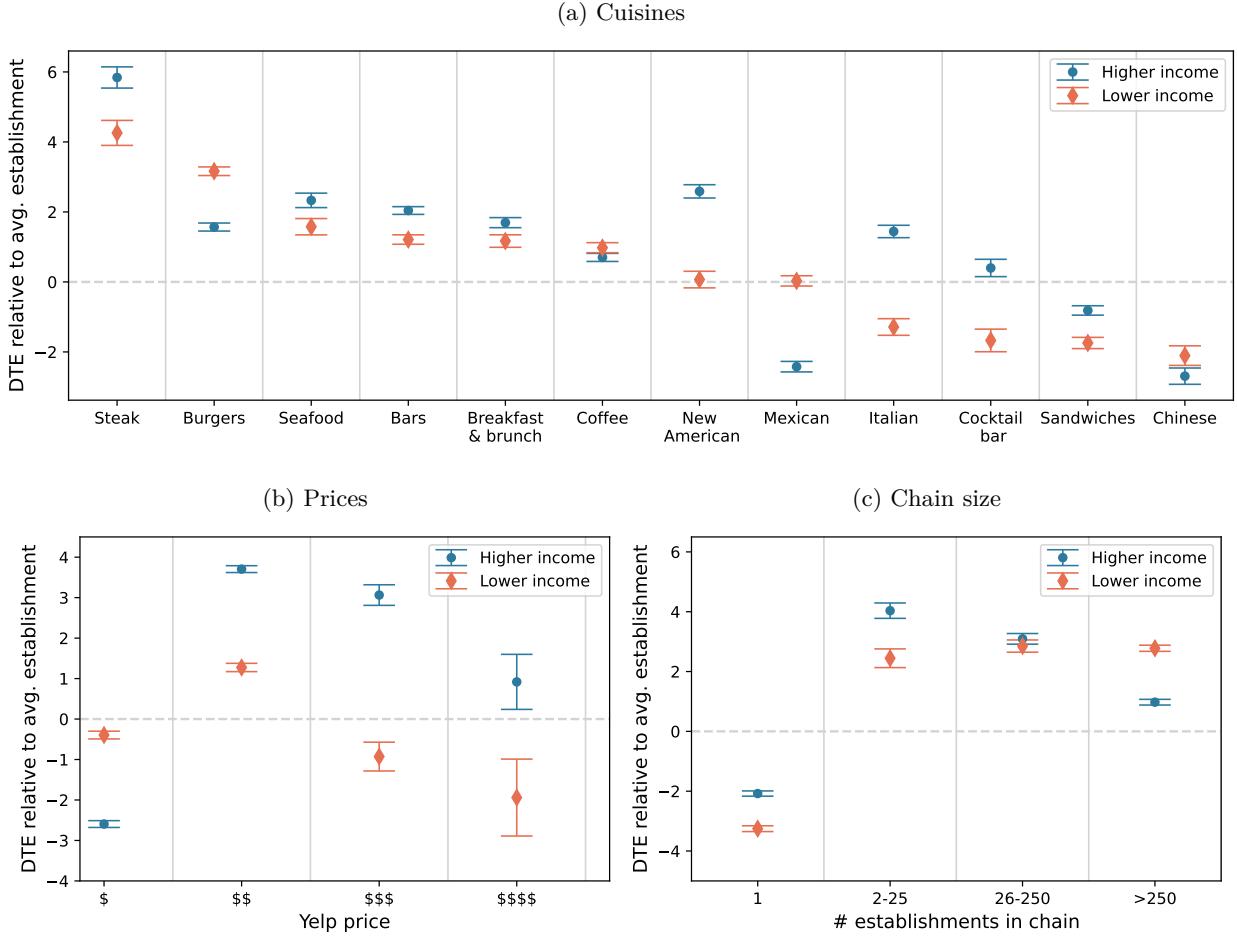
$$\tilde{\gamma}_{jk} \sim \alpha_k X_j + CBSA_j \quad (5.1)$$

where  $X_j$  is always an indicator for whether an establishment has a given characteristic, e.g., whether its cuisine is Mexican or its brand is Trader Joe’s. I run each regression separately for each characteristic and income group. The coefficient  $\alpha_k$  captures the average change in DTEs for an establishment possessing a given characteristic, relative to the average establishment.

Figure 3 shows how preferences for restaurants depend on Yelp characteristics and chain size. Panel (a) shows preferences for different cuisine labels. Higher income residents have relatively stronger preferences for cuisines such as New American, Seafood, and Italian, while lower income residents have relatively stronger preferences for burgers, Mexican, and Chinese food. Panel (b) shows that lower income residents tend to prefer restaurants with price level of \$\$ and dislike pricier restaurants, while higher income residents prefer the pricier \$\$\$-\$\$\$\$ restaurants over cheap, \$ restaurants. These trends echo the findings of Davis et al. (2019), who use data on Yelp reviewers to show that \$\$ restaurants are generally preferred to \$ restaurants and that residents of higher

income Census tracts exhibit more willingness to pay for the higher priced restaurants.<sup>14</sup> Panel (c) documents preferences for restaurants that belong to chains. In general, chain restaurants have broad appeal relative to establishments with only a single location. For the largest chains ( $>250$  locations nationwide), lower income residents tend to have stronger preferences than higher income residents.

Figure 3: Relationship between DTEs and Yelp characteristics (Restaurants)



*Note:* This figure documents the relationship between the establishments' driving time equivalents (DTEs) for each income group and observables. Each marker corresponds to the coefficient from a regression of DTEs for that device type on whether the restaurant has a given characteristic (e.g., Yelp price of '\$\$\$\$', with CBSA fixed effects. The regressions are each run separately and the comparison group is always all restaurants that do not have a given characteristic. Each marker represents a 95% confidence interval. Cuisines are displayed in order of the point estimates for lower income residents.

I next look at the relationship between preferences and an establishment's brand. Certain brands such as Whole Foods and Starbucks have become commonly associated with gentrification in the popular press (see, for example, [Kolomatsky, 2020](#)). Figure 4 plots regression coefficients for three

<sup>14</sup>[Klopack \(2021\)](#) shows that the Yelp price labels correspond closely with the average amount spent at a restaurant.

types of amenities with prominent brands – restaurants, grocery stores, and general merchandise stores. Panel (a) shows preferences for a number of the top restaurant brands. In general, lower income residents place higher value on these chain restaurants than higher income residents. Of the 50 most common restaurant chains, 76% are preferred by lower income residents. The exceptions include sweetgreen, Starbucks, and Chipotle – for each of these chains, higher income residents tend to be more willing to drive further than lower income residents. However, there is also a substantial vertical component of preferences – while lower income residents tend to prefer Taco Bell and higher income residents tend to prefer Chipotle, both types prefer McDonald's to any other chain.

Panel (b) and (c) documents preferences for grocery brands and general merchandise stores. Here, the vertical component of preferences is even more pronounced. While higher income residents tend to place more value on a Trader Joe's and lower income residents have relatively stronger preferences for Save-A-Lot, both types of residents prefer the average Meijer or Kroger to budget options such as Aldi. Even Whole Foods, which is commonly considered a harbinger of gentrification ([Kolomatsky, 2020](#)), is an above average grocery store for lower income residents and preferred to some budget options like Save-A-Lot. The magnitudes here are substantial; individuals prefer Meijer relative to the average supermarket by over 15 DTEs, consistent with previous research finding that households are willing to travel substantial distances to shop at their preferred grocery store ([Allcott et al., 2019](#)). Within general merchandise, everyone prefers either Costco or Walmart to the many brands of dollar stores. While dollar stores have become ubiquitous over the past decade—especially in lower income neighborhoods—these results suggest that residents would generally prefer to shop at a Walmart or Costco, but may frequent dollar stores when they are substantially closer than a larger establishment like Walmart (and, given the low entry costs for a dollar store, they are able to open enough locations to often be the closer option).

Appendix Table C2 reports the top brands for each subcategory as well as the brands with the largest gap between higher and lower income preferences. Many of these match what one might expect – for automobile dealers, for example, Toyota is the top brand for both higher and lower income residents, but the dealers with the biggest gaps in preferences are Audi (H minus L preferences) and Mitsubishi (L minus H preferences). The fitness center brand with the largest difference between higher and lower income preferences is Pure Barre, while the brand with largest difference between lower and higher income preferences is Planet Fitness.<sup>15</sup>

To the extent we can observe similar characteristics of proposed entrants to a neighborhood,

---

<sup>15</sup>[Couture et al. \(2019\)](#) create a ranking of restaurant brands that places more value on brands frequented by residents of richer block groups. Similar to their results, I find that the restaurant brands with the largest gap between higher and lower income preferences tend to be smaller, gourmet chains including Zoë's Kitchen (1st), Cafe Rio (2nd), and First Watch Restaurants (3rd). While they do not report the full ranking, they also find that Zoë's Kitchen is in the top 3.

these results suggest that the residents most likely to value a new entrant is predictable ex-ante.<sup>16</sup> However, for many other categories of amenities—such as barbers and dentists—there are little to no observables available at scale that describe how establishments within the subcategory may differ from each other.

### 5.1.2 Correlation of preferences within each subcategory

Many establishments have few observable characteristics, besides their subcategory. However, I can use the estimated establishment-level preferences for each device type to measure the overall correlation in preferences for establishments within a given subcategory. In subcategories with higher positive correlation in preferences, any given establishment is more likely to have ‘aligned’ preferences between higher and lower income residents.

For each subcategory, Table 2 documents correlation in establishment-level preferences ( $\tilde{\gamma}_{jk}$ ), the estimated standard deviation of the idiosyncratic component of utility, and the percent of each device type’s visits that are to a given subcategory. Each statistic is averaged across CBSAs, weighted by the total number of visits in that subcategory. Subcategories are ordered within a main category by their total number of visits.

First, preferences are positively correlated in all 45 subcategories of amenities. However, the degree of correlation in establishment level preferences for higher and lower income residents varies substantially across subcategories. The correlation is highest in categories dominated by large establishments, such as malls, general merchandise, spectator sports, and gambling. The correlation is lowest for many personal services, such as mechanics, car washes, drycleaning/laundry, and dentists, as well as for book/news stores and museums. Across all categories, the average unweighted correlation in preferences is 0.65. When weighted by the number of visits in a subcategory, the average correlation is 0.82, as larger categories tend to have more correlated preferences.

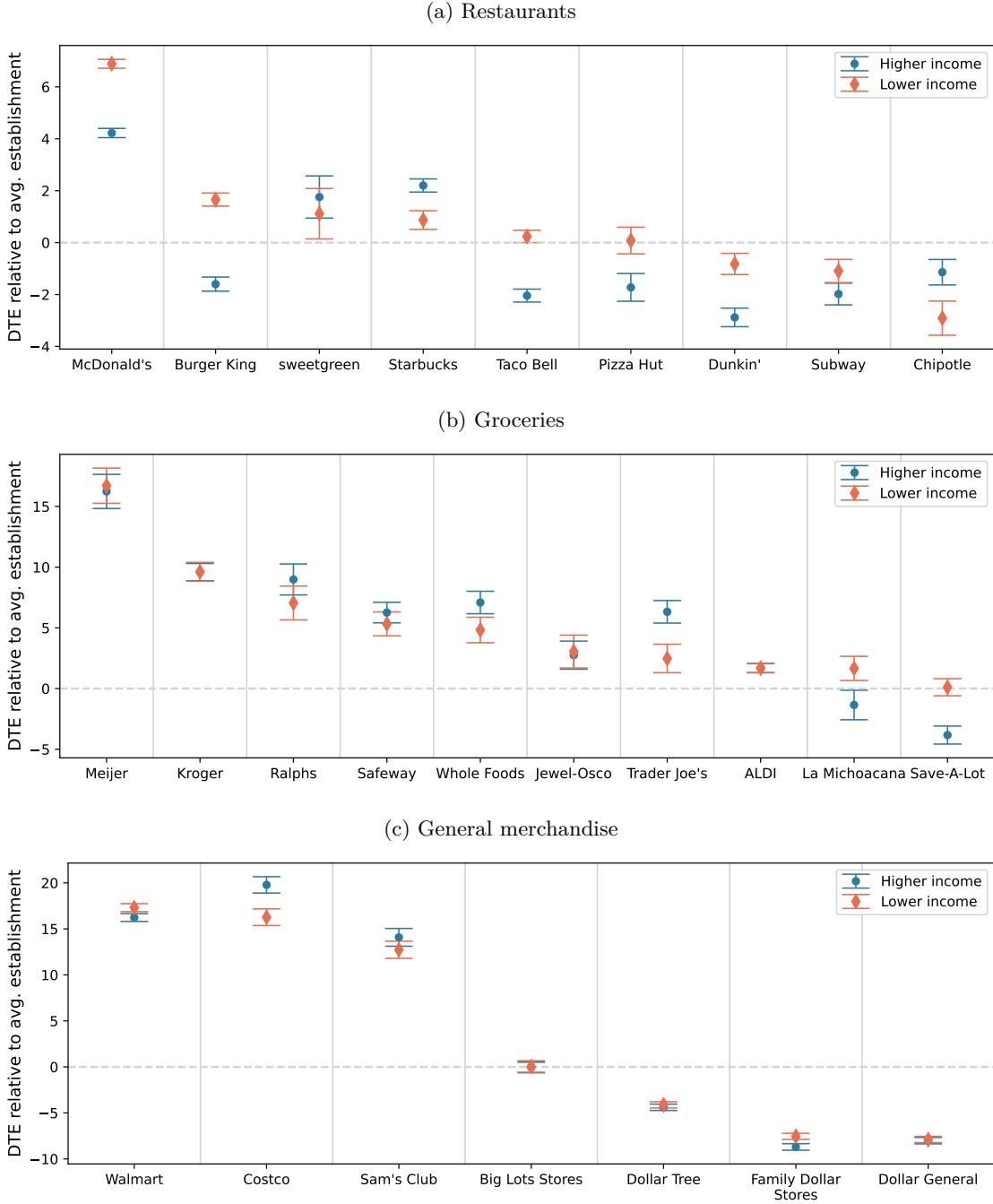
Normalizing utility by the distaste for driving time allows the variance of the idiosyncratic terms to differ across subcategories and across types of individuals. The middle two columns of Table 2 reports the estimated standard deviation of the idiosyncratic component of utility for each subcategory. The idiosyncratic component of utility plays the largest role in many of the entertainment categories, such as spectator sports, gambling, performing arts, and amusement parts.<sup>17</sup> It is a relatively smaller factor for categories such as car washes, libraries, drycleaning/laundry, and

---

<sup>16</sup>Preferences may also be informative about which establishments exit. In Appendix Section C.1, I use a small sample of observed exits to demonstrate how preferences are predictive of which establishment will exit as neighborhood gentrifies.

<sup>17</sup>Agarwal et al. (2017) find that credit card expenditure shares decline more rapidly with distance in sectors transacted more frequently. In the model presented here, the implied visit shares will decline less rapidly with distance if the variance of the idiosyncratic component is larger. Consistent with their finding that that distance matters least for entertainment and sellers of durable goods such as cars, I estimate that the variance of the idiosyncratic components is high in these subcategories.

Figure 4: Relationship between DTEs and brands



*Note:* This figure documents the relationship between establishments' driving time equivalents (DTEs) for different income groups and brands. Each marker corresponds to the coefficient from a regression of DTEs for that device type on whether an establishment is affiliated with a given brand, with CBSA fixed effects. The regressions are each run separately and the comparison group is always all other establishments within that subcategory. Each marker represents a 95% confidence interval. Brands are displayed in order of the point estimates for lower income residents.

gas stations. On average, idiosyncratic preferences play a similar role for higher and lower income residents (weighted average of 9.26 and 9.39, respectively).

Finally, these correlations are for preferences for establishments conditional on consuming within a given subcategory. Individuals may also have heterogeneous preferences across different subcategories of amenities. The final three columns of Table 2 provides some evidence that this is the case by documenting the percent of each device type's visits that are within each subcategory. Higher and lower income residents spend a similar percent of their visits to amenities to various subcategories of restaurants (11.4% and 11.7%, respectively) and of shops (55.4% and 55.8%). Some differences arise for personal services and entertainment places – higher income residents allot 14.0% of their visits to personal services and 19.2% to entertainment, while higher income residents allot 15.9% and 16.6% of their visits to the two categories, respectively. Many of the largest differences across subcategories are within the entertainment category, where higher income residents spend proportionately more time at golf courses, fitness centers, spectator sports, and performing arts than lower income residents. Within restaurants, lower income residents spend slightly more of their visits at limited-service restaurants. For shopping, malls, which includes both monolithic malls filled with many stores as well as smaller strip malls, represent an disproportionately large share of amenity visits for all device types.<sup>18</sup> The different levels of consumption for subcategories of amenities motivates nesting the establishment level estimates into a larger hierarchical model, where preferences can vary across different categories and subcategories.

## 5.2 Who likes which neighborhoods?

### 5.2.1 Correlation in preferences for neighborhoods

The choice model provides a framework for comparing the value of amenity access for residents of different neighborhoods through the NAQI values defined in Equation 4.9.<sup>19</sup> For each income group, NAQI units correspond to weekly minutes of driving time relative to the median neighborhood in a CBSA; a NAQI value of 50 for higher income residents implies that higher income residents value the amenity access of that neighborhood relative to their median neighborhood equivalently to saving 50 minutes of driving time each week.

---

<sup>18</sup>This is in part due to data limitations; for many malls, the POI data do not provide separate polygons for each store within a mall and so visits to any store within the mall are only identified as being at the mall itself, not at the specific store. See Appendix Section A.1.3 for more details on how stays are matched to individual POIs.

<sup>19</sup>Appendix Table C3 documents the average estimated arrival rates ( $\bar{\lambda}$ ) and measures of nest independence ( $\bar{\rho}$ ) across CBSAs. The arrival rate of choice opportunities is only identified up to the normalization of the outside option to zero. As the utility of the outside option is always normalized to zero, the extent to which the arrival rates vary across types of individuals or across times of the week will be determined both by true differences in the frequency of choice opportunities as well as differences in the value of the outside option. The subcategory and category intercepts are not directly interpretable, as they capture both the propensity to consume in a given subcategory as well as the re-scaling needed thanks to the coefficients of integration in the DTEs that link the lower level results to the upper levels (Heiss, 2002).

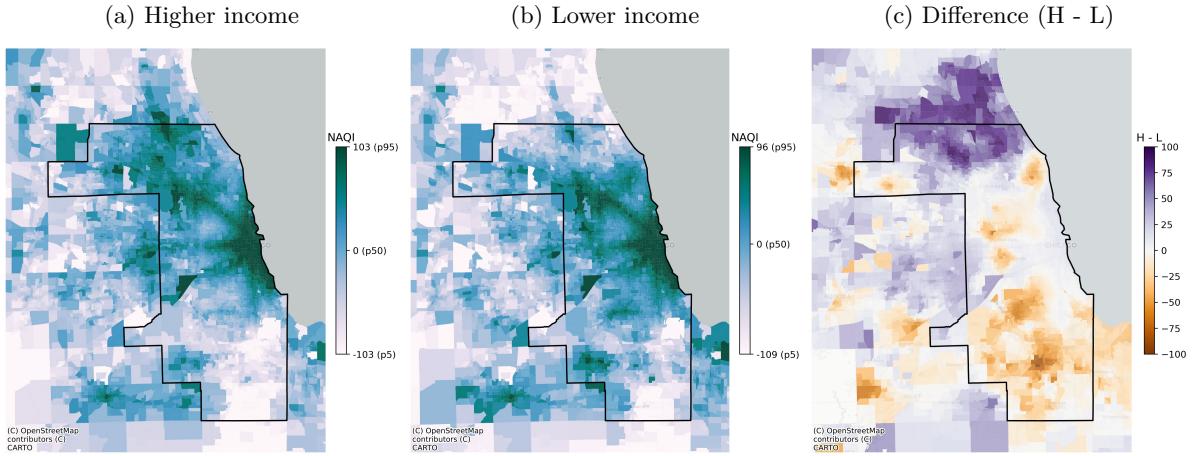
Table 2: Preferences for establishments within a subcategory

	Corr(H, L)	Std. dev. of error term		% of all visits	
		H	L	H%	L%
<b>Restaurants</b>					
Full-service	0.760	10.161	10.545	6.80	6.84
Limited-service	0.692	9.637	9.858	1.95	2.31
Cafes; snacks	0.676	9.685	10.177	1.37	1.30
Drinking places	0.707	11.539	11.374	1.23	1.29
<b>Shops</b>					
Malls	0.935	8.216	8.436	39.90	37.72
General merch., warehouse clubs	0.926	6.692	6.984	2.69	3.60
Gas stations, convenience	0.731	10.851	11.286	1.99	2.96
Grocery stores	0.777	7.436	7.837	1.91	2.16
Clothing, shoes	0.712	10.847	11.111	1.20	1.28
Building materials, gardening	0.826	7.408	7.982	1.25	1.41
Automobile dealers	0.781	12.624	12.659	1.16	1.19
Pharmacies	0.672	7.563	7.852	0.88	1.06
Beauty, glasses, personal care	0.735	8.100	8.137	0.85	0.89
Sporting goods, hobby, music	0.604	10.087	9.939	0.65	0.60
Furniture, furnishings	0.581	10.268	9.884	0.64	0.60
Department	0.797	6.780	6.745	0.63	0.58
Beer, wine, liquor, tobacco	0.483	8.146	7.570	0.36	0.37
Jewelry, luggage, leather	0.599	9.519	9.533	0.35	0.33
Electronics, appliances	0.528	8.637	8.023	0.31	0.37
Used merchandise	0.539	8.191	7.897	0.28	0.38
Office supplies	0.410	8.770	8.244	0.22	0.22
Books, news	0.346	6.288	5.296	0.11	0.11
<b>Personal services</b>					
Hospitals	0.906	12.427	11.747	5.01	6.17
Salons, barbers	0.670	8.501	8.664	2.92	3.33
Religious organizations	0.515	9.427	9.404	2.03	1.95
Banks	0.645	8.367	8.519	1.20	1.46
Other health practitioners	0.530	9.444	9.373	1.05	0.96
Dentists	0.471	8.816	8.389	0.85	0.80
Mechanics	0.297	8.563	8.510	0.39	0.50
Telecom	0.457	7.310	6.272	0.20	0.25
Accounting, taxes, payroll	0.415	6.681	6.052	0.17	0.22
Drycleaning, laundry	0.411	5.742	5.567	0.11	0.13
Car washes	0.247	4.817	4.214	0.09	0.10
<b>Entertainment</b>					
Parks	0.821	9.185	9.437	8.79	8.89
Golf courses	0.840	9.404	10.893	3.89	2.02
Fitness centers	0.741	9.043	8.949	2.68	1.88
Spectator sports	0.844	20.501	17.565	0.77	0.59
Gambling	0.891	16.395	15.344	0.47	0.72
Performing arts	0.824	18.982	18.107	0.56	0.48
Amusement parks	0.821	18.510	17.699	0.45	0.45
Movie theaters	0.770	9.060	8.681	0.56	0.49
Other amusement, recreation	0.619	10.461	9.731	0.39	0.40
Bowling	0.603	9.345	7.686	0.27	0.27
Libraries	0.470	5.851	5.492	0.24	0.25
Museums, zoos, gardens	0.442	11.718	10.524	0.15	0.13

*Note:* This table documents estimation results for each subcategory of amenities. The correlations in preferences and standard deviations of error terms are aggregated across CBSAs by taking the average of the CBSA-level estimation results weighted by the total number of visits to that subcategory in a CBSA. Correlations exclude a small number of establishments that had zero visits by either H or L residents. H and L refer to higher and lower income residents (respectively).

Figure 5 illustrates the NAQI values for the Chicago CBSA, zoomed in on Cook County. The most valuable neighborhoods for both higher and lower income residents are those closer to the downtown area – residents value the amenity access of living downtown by upwards of 100 minutes of saved weekly driving time over the median neighborhood. Block groups along the major arterial highways leading into downtown Chicago are also high value, thanks to the access they provide the amenity-rich urban core. Panel (c) maps the difference between NAQI values for higher and lower income residents. While both residents value the amenities of the urban core, differences in preferences arise as we move away from the core towards less dense neighborhoods. Consistent with the income sorting patterns in Chicago, lower income residents have higher relative values for the amenities of South Chicago and higher income residents have higher values for the amenities in the more Northern parts of the city. Appendix Figure C.5 replicates these maps for the San Francisco Bay Area, Los Angeles, and Boston.

Figure 5: NAQI: Chicago



*Note:* This figure illustrates the estimated block group level NAQI values in Equation 4.9 for the Chicago-Naperville-Elgin, IL-IN-WI CBSA, zoomed in to Cook County (outlined in black). NAQI units correspond to minutes of weekly driving time relative to the median neighborhood in a CBSA.

While Figure 5 shows some systematic differences in preferences for neighborhoods outside of Chicago's urban core, in general preferences for neighborhoods' access to amenities exhibit strong positive correlations; in Chicago, higher and lower income preferences for neighborhoods have a correlation coefficient of 0.90. This level of correlation in NAQI values is representative of other cities – the average within-CBSA correlation between higher and lower income NAQI values is 0.90 with a standard deviation of 0.05 (weighting by CBSA population).

To contextualize the magnitude of the correlation in preferences for neighborhoods, Table 3 documents the average value of moving from the 25th percentile neighborhood in the within-CBSA distribution of neighborhoods for a given type to the 75th percentile in the within-CBSA distribution of either the same income group (along the diagonal) or in the distribution for the

opposite income group (on the off-diagonal). On average, higher income residents value moving from the 25th to 75th percentile in their own distribution equivalently to 78.2 minutes of weekly driving time and to the 75th percentile neighborhood of the lower income distribution equivalently to 73.3 weekly minutes. Lower income residents value moving from the 25th to 75th percentile neighborhood equivalently to 71.7 minutes and value moving to the 75th percentile neighborhood of higher income residents equivalently to 67.0 minutes. Moving to the opposite device income type's 25th percentile neighborhood is only 6-7% worse than moving to the 25th percentile neighborhood of given device type's distribution of neighborhood values. However, these results do not preclude higher income residents placing more absolute weight on the amenity value of a neighborhood when choosing where to live – given that higher income residents visit 10-30% more amenities of each category in a quarter, it would be sensible that they value amenity access more. In the context of earlier findings on the role of amenities in the the return of young, wealthy residents to urban cores (Couture et al., 2019), the level of correlation observed here implies these trends are likely due to higher income residents placing increasing weight on amenity access in recent decades when choosing where to live, which may price out lower income residents who place similar *relative* value on the amenity access of the urban core but are unwilling to pay as much for amenity access (Couture and Handbury, 2020).

Table 3: Value of moving from p25 to p75 neighborhood

Device type / p25 distribution	Distribution of NAQI values for p75	
	Higher income	Lower income
Higher income	78.22 [34.06]	73.28 [34.27]
Lower income	67.03 [46.38]	71.66 [44.97]

*Note:* This table documents the average value of moving from the 25th percentile neighborhood in a given within-CBSA distribution of NAQI values to 75th percentile neighborhood in the opposite income group's distribution of NAQI values. Units correspond to the value of a minute of driving time. The results are computed within-CBSA and then averaged across CBSAs using total population as weight. The weighted across-CBSA standard deviations are in brackets.

### 5.2.2 Characteristics of neighborhoods with high value amenity access

What types of neighborhoods tend to have good access to amenities? I examine the relationship between amenity access for each income group and neighborhood characteristics by running a series of regressions of NAQI values on characteristics such as density, median household income, and rent. Each regression uses data from all 100 CBSAs and includes CBSA fixed effects. The values for non-logged variables correspond to a one standard deviation increase in the covariate.

The first three columns of Table 4 document the results for univariate regressions. Much as Figure 5 demonstrated for Chicago, the two features most predictive of high quality amenity access are population density and the distance to city hall. A doubling population density is

associated with a  $\sim$ 30 minute increase in the NAQI value of a neighborhood for each income group. Regressions with either of these covariates have  $R^2$  values around 0.35. For both higher and lower income residents, the value of a neighborhood's amenity access is also increasing in rent and education and decreasing in household income, fraction white-alone, and median age. Similar to the finding in Figure 1 that higher income residents have fewer establishments nearby, Table 4 documents that the value of a neighborhood's amenities is lower in higher income neighborhoods. The third column shows that the absolute gap between higher and lower income residents—a measure of ‘disagreement’ over neighborhood quality—tends to be lower in dense neighborhoods that are closer to city hall and higher in more white, younger neighborhoods.

Many of the covariates are correlated with density. The latter three columns control for log density and and log distance to city hall. With these controls, neighborhoods with higher income, rent, college graduates, and more white or older residents tend to have higher quality access to amenities. Despite the levels of positive correlation in NAQI values, some systematic differences do emerge. Once controlling for density, the relationship between NAQI values and neighborhood income, rent, fraction white, and median age is about twice as strong for higher income residents than lower income residents. For example, each standard deviation increase in neighborhood income is associated with a 7.9 minute increase in NAQI value for higher income residents versus 4.3 minutes for lower income residents.

### 5.3 Counterfactual: tailoring amenities to higher income residents

A common concern with gentrifying neighborhoods is that the tailoring of amenities to higher income entrants will amplify the overall welfare inequality in the neighborhood by crowding out amenities enjoyed by the lower income incumbents. To evaluate the potential magnitude of welfare effects through this channel, I simulate a counterfactual world in which the top 25% of establishments for lower income residents within each subcategory are replaced with duplicates of the top 25% of establishments for higher income residents. The new establishments are set in the same location in the CBSA as the ones they replace, but I switch out the establishment-level preference parameters for each income group with the parameters from the new establishments. Subcategories with highly correlated preferences will have substantial overlap in the top 25% of establishments for higher and lower income residents.

I find that this fairly extreme version of tailoring has a modest effect on the overall utility residents derive from neighborhood amenities. In the average CBSA, a higher income resident’s utility increases by an amount corresponding to 18.1 minutes of weekly driving time and a lower income resident’s utility decreases by 21.4 minutes of weekly driving time. If residents have an hourly value of time of \$20, the changes in utility correspond to +\$313 per year for higher income residents and  $-\$371$  per year for lower income residents.<sup>20</sup> In a similar exercise, Davis et al. (2019) simu-

---

<sup>20</sup>Goldszman et al. (2020) estimate a value of time of \$19 using variation in the wait times for Lyft rides.

Table 4: Relationship between NAQI values and neighborhood characteristics

Covariate	Univariate			Controls: density & city hall distance		
	H	L	H - L	H	L	H - L
Log population density	29.5364 (SE: 0.174) [R2: 0.3481]	30.286 (SE: 0.1729) [R2: 0.3793]	-3.5268 (SE: 0.0771) [R2: 0.3233]			
Log miles to city hall	-58.105 (SE: 0.2703) [R2: 0.3726]	-55.9451 (SE: 0.269) [R2: 0.3607]	4.737 (SE: 0.1265) [R2: 0.3108]			
Median income	-4.311 (SE: 0.1954) [R2: 0.0269]	-7.9668 (SE: 0.1868) [R2: 0.0289]	0.2498 (SE: 0.0699) [R2: 0.2946]	7.9316 (SE: 0.1698) [R2: 0.4565]	4.3337 (SE: 0.16) [R2: 0.4653]	-1.2305 (SE: 0.0777) [R2: 0.322]
Median rent	4.5393 (SE: 0.2747) [R2: 0.029]	1.7062 (SE: 0.262) [R2: 0.021]	-0.424 (SE: 0.093) [R2: 0.3133]	7.2876 (SE: 0.2155) [R2: 0.4444]	4.4207 (SE: 0.2001) [R2: 0.4471]	-0.7009 (SE: 0.0926) [R2: 0.3367]
Frac. college grad	12.1012 (SE: 0.1951) [R2: 0.0476]	7.0334 (SE: 0.1905) [R2: 0.0269]	-0.1759 (SE: 0.0716) [R2: 0.2982]	10.6795 (SE: 0.1481) [R2: 0.465]	5.9041 (SE: 0.1425) [R2: 0.4656]	-0.2305 (SE: 0.0718) [R2: 0.3239]
Frac. white-alone	-14.9958 (SE: 0.2041) [R2: 0.0586]	-17.9177 (SE: 0.1963) [R2: 0.0714]	1.8314 (SE: 0.0757) [R2: 0.3011]	7.0405 (SE: 0.1755) [R2: 0.4527]	3.4586 (SE: 0.1705) [R2: 0.4611]	-0.3523 (SE: 0.0835) [R2: 0.3239]
Median age	-12.0876 (SE: 0.2199) [R2: 0.0488]	-12.9717 (SE: 0.2133) [R2: 0.0492]	1.2898 (SE: 0.0806) [R2: 0.2998]	3.4936 (SE: 0.1848) [R2: 0.4481]	2.8647 (SE: 0.1764) [R2: 0.4607]	-0.5525 (SE: 0.083) [R2: 0.3241]

*Note:* This table documents results from a series of regressions of the block group level NAQI values on characteristics of the block group. Population and demographics data are from the 2019 5-year ACS. The distance to city hall is crow-flies, based on the city hall location reported on Google Maps for the largest city in each CBSA. Regressions include fixed effects for the CBSA. For non-log covariates, the values correspond to a one standard deviation increase.

late replacing both restaurants and residents of neighborhoods surrounding majority-Black Harlem with those from majority-white Upper East Side and find that the change in restaurant prices and cuisines has only a small effect on welfare relative to the social frictions that arise with the shift of surrounding residents from mostly Black to mostly white. Similar to their findings, the results presented here suggest that the direct effect of amenities tailoring to higher income residents on spatial inequality is relatively small thanks to the correlation in preferences for amenities; however, these partial equilibrium effects may be further amplified (or dampened) by general equilibrium channels such as rent increases and residential re-sorting.

#### 5.4 Robustness tests

I conduct a number of robustness tests that either vary the definition of income and the sample of trips used for estimation. For each each test, I use a subset of ten CBSAs sampled randomly from each decile of CBSA population to limit the computational burden of re-estimating the entire set of

preferences.<sup>21</sup> In Table 5, I report the correlation in establishment-level preferences, correlation in NAQI values, and the results from the counterfactual tailoring to higher income residents. The first row documents the baseline results for these ten CBSAs, which are similar to the results presented for the entire sample.

#### 5.4.1 Alternative income groups

The first set of robustness tests varies the definition of higher and lower income. A concern with using imputed income is that the correlation in estimated preferences will be sensitive to the level of measurement error. At the extreme, if I were to randomly assign income labels to devices then I would find that preferences are perfectly correlated. To help address this concern, I use three alternative definitions of higher and lower income. First, I forego the parcel-level income estimates and instead use whether a device's block group median income is above or below the median income for the CBSA. Second, I again use the block group median household income, but exclude devices living in mixed-income block groups, which I define as having less than two-thirds of either income group. Finally, I return to using parcel-level income estimates, but define higher and lower income based on being in the top or bottom tercile of income and exclude devices in the middle tercile.<sup>22</sup>

In each case, the story is similar: preferences for establishments are less correlated than preferences for neighborhoods and tailoring establishments to higher income residents has a modest effect on overall welfare. When using the tails of the income distribution, the correlation in preferences falls and the effects on welfare are amplified. The upper tercile of the income distribution values the increased tailoring equivalent to 32.3 minutes of weekly driving time (+\$559 per year) while the welfare of residents in the bottom tercile declines by the equivalent of 50.71 minutes of weekly driving time (-\$878 per year). As we would expect, 'higher income' and 'lower income' residents are not a monolith and preferences at the tails of the income distribution are more polarized than the preferences of residents in the middle tercile.

#### 5.4.2 Subsets of amenity visits

The second set of tests uses specific subsets of visits to amenities. First, I re-estimate the model using only the first visit to an establishment by each device (57% of all visits). The model assumes that each choice is made independent of other choices an individual has previously made; however, an individual's the idiosyncratic draws for a given location are likely correlated across time. In a similar model, [Davis et al. \(2019\)](#) show that serial correlation in the idiosyncratic errors can lead to attenuation bias in the coefficient estimates. In practice, subsetting to only the first trip by each device has little effect on the bottom line results.

<sup>21</sup>The ten CBSAs are, in order of population, Chicago-Naperville-Elgin, Seattle-Tacoma-Bellevue, Kansas City, Cleveland-Elyria, New Orleans-Metairie, Albuquerque, Knoxville, Cape Coral-Fort Myers, Madison, and Durham-Chapel Hill.

<sup>22</sup>Appendix Section [A.2.4](#) provides details on how devices are split into income terciles.

Second, I subset to visits where both the previous and subsequent stops are at home (34% of all visits) and set the disutility of driving time from work to zero. Individuals may combine multiple stops into a single trip to amortize the disutility of travel time. Treating each stop as an independent trip would bias upwards the estimated value of establishment that is frequently visited as part of a chain of trips (e.g., an ice cream shop next to a group of retail stores). Using only trips that start and end at home reduces both the correlation in preferences and the effect on welfare of tailoring establishments to higher income devices. However, these differences appear to be primarily driven by the small sample size introducing additional measurement error in establishment-level preferences. When using a random third of trips (the final row), the results are similar to when using only the trips that start and end at home. In Appendix Section A.1.4, I discuss the difficulties with identifying chains of trips to multiple establishments in GPS data.

Table 5: Robustness tests

	Establishments corr.		NAQI corr. All	Tailoring counterfactual	
	Unweighted	Weighted		Lower income	Higher income
<b>Baseline</b>	0.634	0.805	0.914	19.19	-24.26
<b>Alternative income measures</b>					
Block group income	0.645	0.803	0.899	19.33	-30.34
Block group income, excluding mixed-income block groups	0.403	0.612	0.855	32.27	-50.71
Top and bottom income terciles	0.470	0.674	0.843	23.26	-44.45
<b>Subset of visits</b>					
First trip to establishment by each device	0.621	0.787	0.861	22.10	-17.85
Trips that start and end at home	0.480	0.730	0.924	7.69	-8.94
Random subset of one-third of trips	0.407	0.709	0.840	11.82	-10.66

*Note:* This table documents the results of a battery of robustness tests. The tailoring counterfactual involves replacing the top 25% of establishments for lower income residents in each subcategory with replicas of the top 25% of establishments for higher income residents. The units are minutes of weekly driving time.

## 6 Conclusion

This paper studies heterogeneity in preferences for neighborhood amenities at two levels: 1) preferences for individual establishments and 2) preferences for neighborhoods' overall access to amenities. At the establishment-level, I find that higher and lower income residents have positively correlated preferences for all 45 different subcategories of amenities; however, preferences are more correlated within some types, such as malls and general merchandise stores, than others, such as many per-

sonal services. For amenities where preferences are less positively correlated, there is more scope for new establishments in gentrifying neighborhoods to tailor to the preferences of the newer, higher income residents to the detriment of the incumbent, lower income residents.

At the neighborhood-level, preferences exhibit strong positive correlation. While higher and lower income residents may disagree about which restaurants they prefer, for example, they tend to agree on which neighborhoods offer the best access to amenities. Neighborhoods with high quality amenity access tend to be denser, closer to the urban core, younger, higher education, and with higher rents.

There are two key margins that are not addressed in this paper: 1) the supply-side decisions of local amenities and 2) the effect of local amenities on housing markets and residential choices. Researchers interested in the latter question may find value in the estimated NAQI values, which summarizes the quality of each neighborhood's access to amenities. For example, these can serve as an input into models of residential choice to study heterogeneity in preferences for living in neighborhoods with better access to amenities.

## References

- Abadi, Martin, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng,** “TensorFlow: A system for large-scale machine learning,” *12th USENIX symposium on operating systems design and implementation*, 2016, p. 21.
- Abbiasov, Timur and Dmitry Sedov,** “Do local businesses benefit from sports facilities? The case of major league sports stadiums and arenas,” *Regional Science and Urban Economics*, 2022, p. 103853.
- Agarwal, Sumit, J Bradford Jensen, and Ferdinando Monte,** “Consumer mobility and the local structure of consumption industries,” Technical Report, National Bureau of Economic Research 2017.
- Ahlfeldt, Gabriel M, Stephen J Redding, Daniel M Sturm, and Nikolaus Wolf,** “The economics of density: Evidence from the Berlin Wall,” *Econometrica*, 2015, 83 (6), 2127–2189.
- Albouy, David,** “Are big cities bad places to live? Estimating quality of life across metropolitan areas,” Technical Report, National Bureau of Economic Research 2008.
- Allcott, Hunt, Levi Boxell, Jacob Conway, Billy Ferguson, Matthew Gentzkow, and Benny Goldman,** “What Explains Temporal and Geographic Variation in the Early US Coronavirus Pandemic?,” Technical Report w27965, National Bureau of Economic Research, Cambridge, MA October 2020.
- , **Rebecca Diamond, Jean-Pierre Dubé, Jessie Handbury, Ilya Rahkovsky, and Molly Schnell,** “Food Deserts and the Causes of Nutritional Inequality\*,” *The Quarterly Journal of Economics*, November 2019, 134 (4), 1793–1844.
- Almagro, Milena and Tomas Dominguez-Iino,** “Location Sorting and Endogenous Amenities: Evidence from Amsterdam,” *Working paper*, 2022.
- Athey, Susan, Billy Ferguson, Matthew Gentzkow, and Tobias Schmidt,** “Estimating experienced racial segregation in US cities using large-scale GPS data,” *Proceedings of the National Academy of Sciences*, 2021, 118 (46).
- , **David Blei, Robert Donnelly, Francisco Ruiz, and Tobias Schmidt,** “Estimating Heterogeneous Consumer Preferences for Restaurants and Travel Time Using Mobile Location Data,” *AEA Papers and Proceedings*, May 2018, 108, 64–67.
- Atkin, David, Keith Chen, and Anton Popov,** “The Returns to Serendipity: Knowledge Spillovers in Silicon Valley,” *Working paper*, 2020.
- Autor, David, Christopher Palmer, and Parag Pathak,** “Gentrification and the Amenity Value of Crime Reductions: Evidence from Rent Deregulation,” Technical Report w23914, National Bureau of Economic Research, Cambridge, MA October 2017.

- Baum-Snow, Nathaniel and Daniel Hartley**, “Accounting for central neighborhood change, 1980–2010,” *Journal of Urban Economics*, 2020, 117, 103228.
- Ben-Akiva, Moshe E**, “Structure of passenger travel demand models.” PhD dissertation, Massachusetts Institute of Technology 1973.
- Birant, Derya and Alp Kut**, “ST-DBSCAN: An algorithm for clustering spatial-temporal data,” *Data & Knowledge Engineering*, January 2007, 60 (1), 208–221.
- Brueckner, Jan K., Jacques-François Thisse, and Yves Zenou**, “Why is central Paris rich and downtown Detroit poor?,” *European Economic Review*, January 1999, 43 (1), 91–107.
- Brummet, Quentin and Davin Reed**, “The effects of gentrification on the well-being and opportunity of original resident adults and children,” *FRB of Philadelphia Working Paper*, 2019.
- Caetano, Gregorio and Vikram Maheshri**, “Gender segregation within neighborhoods,” *Regional Science and Urban Economics*, 2019, 77, 253–263.
- Chang, Serina, Emma Pierson, Pang Wei Koh, Jaline Gerardin, Beth Redbird, David Grusky, and Jure Leskovec**, “Mobility network models of COVID-19 explain inequities and inform reopening,” *Nature*, November 2020.
- Chen, M Keith, Kareem Haggag, Devin G Pope, and Ryne Rohla**, “Racial disparities in voting wait times: evidence from smartphone data,” *Review of Economics and Statistics*, 2022, 104 (6), 1341–1350.
- Chen, Tianqi and Carlos Guestrin**, “Xgboost: A scalable tree boosting system,” in “Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining” 2016, pp. 785–794.
- Cook, Cody, Lindsey Currier, and Edward L Glaeser**, “Urban mobility and the experienced isolation of students and adults,” Technical Report, National Bureau of Economic Research 2022.
- Couture, Victor**, “Valuing the Consumption Benefits of Urban Density,” *Working Paper*, 2016.
- and Jessie Handbury, “Urban revival in America,” *Journal of Urban Economics*, 2020, 119, 103267.
- , Cecile Gaubert, Jessie Handbury, and Erik Hurst, “Income Growth and the Distributional Effects of Urban Spatial Sorting,” Technical Report w26142, National Bureau of Economic Research, Cambridge, MA August 2019.
- , Jonathan I. Dingel, Allison Green, Jessie Handbury, and Kevin R. Williams, “JUE Insight: Measuring movement and social contact with smartphone data: a real-time application to COVID-19,” *Journal of Urban Economics*, February 2021, p. 103328.
- Davis, Donald R., Jonathan I. Dingel, Joan Monras, and Eduardo Morales**, “How Segregated Is Urban Consumption?,” *Journal of Political Economy*, August 2019, 127 (4), 1684–1738.
- Diamond, Rebecca**, “The Determinants and Welfare Implications of US Workers’ Diverging Location Choices by Skill: 1980–2000,” *American Economic Review*, March 2016, 106 (3), 479–

524.

- Ellen, Ingrid Gould and Katherine O'Regan**, "Crime and urban flight revisited: The effect of the 1990s drop in crime on cities," *Journal of Urban Economics*, November 2010, 68 (3), 247–259.
- Fogli, Alessandra and Veronica Guerrieri**, "The End of the American Dream? Inequality and Segregation in US Cities," Technical Report w26143, National Bureau of Economic Research, Cambridge, MA August 2019.
- Glaeser, Edward L., Hyunjin Kim, and Michael Luca**, "Nowcasting Gentrification: Using Yelp Data to Quantify Neighborhood Change," *AEA Papers and Proceedings*, May 2018, 108, 77–82.
- , **Jed Kolko, and Albert Saiz**, "Consumer city," *Journal of Economic Geography*, January 2001, 1 (1), 27–50.
- Glaeser, Edward, Michael Luca, and Erica Moszkowski**, "Gentrification and Neighborhood Change: Evidence from Yelp," Technical Report w28271, National Bureau of Economic Research, Cambridge, MA December 2020.
- Glass, Ruth**, "Aspects of change," in "The gentrification debates," Routledge, 1964, pp. 19–29.
- Goldszmidt, Ariel, John A List, Robert D Metcalfe, Ian Muir, V Kerry Smith, and Jenny Wang**, "The Value of Time in the United States: Estimates from Nationwide Natural Field Experiments," Technical Report, National Bureau of Economic Research 2020.
- Gordinier, Jeff**, "South Bronx Gets High-End Coffee; Is Gentrification Next?," *New York Times*, Mar 2016.
- Guerrieri, Veronica, Daniel Hartley, and Erik Hurst**, "Endogenous gentrification and housing price dynamics," *Journal of Public Economics*, April 2013, 100, 45–60.
- Gupta, Arpit, Stijn Van Nieuwerburgh, and Constantine Kontokosta**, "Take the Q train: Value capture of public infrastructure projects," *Journal of Urban Economics*, 2022, 129, 103422.
- Handbury, Jessie**, "Are Poor Cities Cheap for Everyone? Non-Homotheticity and the Cost of Living Across US Cities," *Econometrica*, 2021, 89 (6), 2679–2715.
- and **David E Weinstein**, "Goods prices and availability in cities," *The Review of Economic Studies*, 2015, 82 (1), 258–296.
- Heiss, Florian**, "Structural Choice Analysis with Nested Logit Models," *The Stata Journal: Promoting communications on statistics and Stata*, September 2002, 2 (3), 227–252.
- Hoelzlein, Mattias**, "Two-sided sorting and spatial inequality in cities," *Working Paper*, 2019.
- Kingma, Diederik P. and Jimmy Ba**, "Adam: A Method for Stochastic Optimization," *arXiv:1412.6980 [cs]*, January 2017. arXiv: 1412.6980.
- Klopack, Ben**, "One size fits all? The value of standardized retail chains," Technical Report, Stanford Working Paper 2021.
- Kolomatsky, Michael**, "Whole Foods and Trader Joe's Downstairs, Higher Rent Upstairs," *New York Times*, Sept 2020.

- Lee, Sanghoon and Jeffrey Lin**, “Natural Amenities, Neighbourhood Dynamics, and Persistence in the Spatial Distribution of Income,” *The Review of Economic Studies*, January 2018, 85 (1), 663–694.
- Lewis, Greg, Bora Ozaltun, and Georgios Zervas**, “Maximum Likelihood Estimation of Differentiated Products Demand Systems,” *arXiv preprint arXiv:2111.12397*, 2021.
- Loshchilov, Ilya and Frank Hutter**, “Decoupled Weight Decay Regularization,” *arXiv:1711.05101 [cs, math]*, January 2019. arXiv: 1711.05101.
- McFadden, Daniel et al.**, “Conditional logit analysis of qualitative choice behavior,” *Analysis of Qualitative Choice Behavior*, 1973.
- McKinnish, Terra, Randall Walsh, and T. Kirk White**, “Who gentrifies low-income neighborhoods?,” *Journal of Urban Economics*, March 2010, 67 (2), 180–193.
- Mikolov, Tomas, Kai Chen, Greg Corrado, and Jeffrey Dean**, “Efficient estimation of word representations in vector space,” *arXiv preprint arXiv:1301.3781*, 2013.
- Miyauchi, Yuhei, Kentaro Nakajima, and Stephen Redding**, “Consumption Access and the Spatial Concentration of Economic Activity: Evidence from Smartphone Data,” Technical Report w28497, National Bureau of Economic Research, Cambridge, MA February 2021.
- Paszke, Adam, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer**, “Automatic differentiation in PyTorch,” *NIPS 2017 Workshop*, 2017, p. 4.
- Raleigh, Helen**, “Gentrification Provokes a Coffee Clash in Denver’s Five Points,” *Wall Street Journal*, Dec 2017.
- Reardon, Sean F., Kendra Bischoff, Ann Owens, and Joseph B. Townsend**, “Has Income Segregation Really Increased? Bias and Bias Correction in Sample-Based Segregation Estimates,” *Demography*, December 2018, 55 (6), 2129–2160.
- Redding, Stephen J and Esteban Rossi-Hansberg**, “Quantitative spatial economics,” *Annual Review of Economics*, 2017, 9, 21–58.
- Relihan, Lindsay**, “Is online retail killing coffee shops? estimating the winners and losers of on-line retail using customer transaction microdata,” *London School of Economics CEP Discussion Paper*, 2022.
- Roback, Jennifer**, “Wages, rents, and the quality of life,” *Journal of Political Economy*, 1982, 90 (6), 1257–1278.
- Rosen, Sherwin**, “Wage-based indexes of urban quality of life,” *Current issues in urban economics*, 1979, pp. 74–104.
- Slickdeals**, “Survey of iPhone and Android owners,” Oct 2018.
- Smith, Rosa Inocencio**, “When a Grocery Store Means Gentrification,” *The Atlantic*, Aug 2016.
- Su, Yichen**, “Measuring the value of urban consumption amenities: A time-use approach,” *Journal of Urban Economics*, 2022, p. 103495.

- Ukueberuwa, Mene**, “Gentrification Is America’s Best Hope,” *Wall Street Journal*, Aug 2020.
- Vigdor, Jacob, Douglas Massey, and Alvin Rivlin**, “Does Gentrification Harm the Poor?,” *Brookings-Wharton Papers on Urban Affairs*, 2002, pp. 133–182.
- Waldfogel, Joel**, “Preference Externalities: An Empirical Study of Who Benefits Whom in Differentiated-Product Markets,” *The RAND Journal of Economics*, 2003, 34 (3), 557.
- , “Who benefits whom in the neighborhood? Demographics and retail product geography,” in “Agglomeration economics,” University of Chicago Press, 2010, pp. 181–209.

# Appendices

<b>A Data appendix</b>	<b>34</b>
A.1 GPS location data . . . . .	34
A.1.1 Pre-processing: from pings to stays . . . . .	34
A.1.2 Assigning home and work locations . . . . .	35
A.1.3 Matching stays to location details . . . . .	35
A.1.4 Trip chains . . . . .	36
A.2 Estimating individual income . . . . .	38
A.2.1 Matching devices to houses . . . . .	38
A.2.2 Estimating market value of homes . . . . .	39
A.2.3 Estimating whether above median income using home characteristics . . . . .	41
A.2.4 Estimating a continuous measure of income . . . . .	44
A.3 Evaluating sample quality . . . . .	45
<b>B Model appendix</b>	<b>45</b>
B.1 Overview of tensor-based estimation frameworks . . . . .	45
B.2 Estimation of lower level of model . . . . .	48
B.3 Estimation of upper levels of model . . . . .	48
<b>C Additional results</b>	<b>50</b>
C.1 Preferences, exits, and gentrification . . . . .	50
C.2 Tourist preferences . . . . .	52
C.3 Additional tables and figures . . . . .	52

## A Data appendix

### A.1 GPS location data

#### A.1.1 Pre-processing: from pings to stays

A raw GPS trace consists of many ‘pings’ attached to a single device. Pings can be triggered either by a user opening specific apps or via apps that share location data in the background. A year of raw data includes over 700 billion pings. To reduce the size of the data while retaining the core information, individual pings that are close together in time and space are collapsed into ‘stays’ by Replica. A stay is composed of a unique device ID, enter time, exit time, latitude, and longitude. For example, a visit to the supermarket may result in a hundred pings, all of which would become a single stay.

To construct stays at locations, pings are combined using an algorithm based on the ST-DBSCAN clustering algorithm, which can discover arbitrarily shaped clusters of successive GPS pings for a given user (Birant and Kut, 2007). The coordinates of a stay are the centroid of all pings assigned to the same cluster by ST-DBSCAN. Pings that occur in transit between stays are discarded.

### A.1.2 Assigning home and work locations

For each day a device appears in the data, I use stays from the 28 days leading up to a given date to assign a most likely home and—where applicable—work location as of that date.

Home assignment consists of three steps. First, I count all overnight (12-5am) stays by a device in an H3 bound of resolution 9. H3 is a hierarchical spatial index, similar to geohashes. Bounds of resolution 9 covers approximately 100 square meters. Second, I label the most common H3 bound as the home bound, so long as it has at least 5 overnight dates. Finally, I identify the exact coordinates of the home by taking the centroid of all overnight stays within the home bound.

Assigning workplaces is more complicated, as work shifts can occur at all hours of the day and, for many occupations, may not involve a single location (e.g., plumbers, taxi drivers, and appraisers). I assign works according to similar heuristics as homes: the *non-home* H3 bound with the most daytime (6am-8pm) stays and average stay length of over 2 hours is labeled the work bound of a device, and the coordinates are the centroid of all stays in this bound.

To determine the home and work locations at the device-quarter level, I use the most commonly assigned home and work for all days for which a device had an assigned home and/or work in a quarter. I exclude any device-quarters for which there is no reliable home location. Of device-quarters with a home assigned, 67% are assigned a work location.<sup>23</sup>

### A.1.3 Matching stays to location details

Stays are matched to details on the exact location—where available—using Point of Interest (POI) data from SafeGraph. SafeGraph aggregates, standardizes, and validates data from many vendors and open source maps on different POIs, which include consumer-facing businesses, schools, parks, hospitals, places of worship, and more.

Each location is also associated with a polygon of its footprint. I use these footprints to match stays to POIs, matching on whether the stay's coordinates are within the polygon during the establishment's open hours. These polygons are not always disjoint, so a single stay may match to multiple locations. There are three primary cases when polygons may overlap. First, when a location is entirely enclosed in another and does not have its own polygon, such as a store in a large mall. In this case, I assign the ‘parent’ location (i.e. the mall) rather than attempting to identify

---

<sup>23</sup>This method for assigning works has drawbacks and cannot identify workplaces for occupations that do not involve a static location (e.g., Uber drivers and mail carriers) or for workers that work night-shifts.

the individual location. Second, stores within a ‘parent’ polygon may have their own individual polygons in the data (e.g., a shopping center where stores are more easily separated). In this case, the staypoint will match to both the parent as well as the individual polygon—since these are overlapping polygons—but the data are sufficient to discern which store within the parent location was visited, so I assign the individual location. Finally, polygons may overlap due to being stacked in multi-story buildings or simply due to noise in the polygon data. In this case, I randomly assign a polygon among the matching set. Of the stays matched to at least one Safegraph POI, 74.3% either match to a single location or a single non-parent location, 9.0% are matched to a parent location, and the remaining 16.7% are randomly assigned a polygon from the set of matches. I exclude visits to POIs that are at a device’s home or work location.

Each location is classified into mutually exclusive subcategories of amenities. I use SafeGraph’s assigned NAICS code to classify each POI in the data. One exception is malls, which do not have a corresponding NAICS code, but are the ‘parent’ POI for many smaller establishments – I use SafeGraph’s ‘Lessors of Real Estate’ category to identify Malls. Finally, I drop establishments with fewer than 5 observed visits. Table A1 documents the subcategories, NAICS codes, and the number of POIs.

#### A.1.4 Trip chains

In modeling amenity consumption, I consider each choice to visit an amenity in isolation. In reality, individuals may combine many stays at different POIs into a single ‘trip,’ allowing them to amortize the cost of driving across multiple stops. This simplification could lead to over-estimates of the value of locations that are often combined with other stops. In this section, I examine the frequency with which trip chains are observed. Even for devices with fairly reliable coverage, I find that the vast majority of trips consist of only a single observed stop at a POI.

I combine stays at amenities into ‘trips’ by labeling a stay to be the start of a new trip if 1) the stay is at a POI that is not home or work and 2) the previous stay was at either home or work or the time between stays was over two hours. This final restriction is to exclude potential trips with large gaps between stays, which is more likely to be due to issues with GPS coverage than long travel times between locations. This definition identifies 165.0 million trips with an average length of 1.10 stays at POIs.

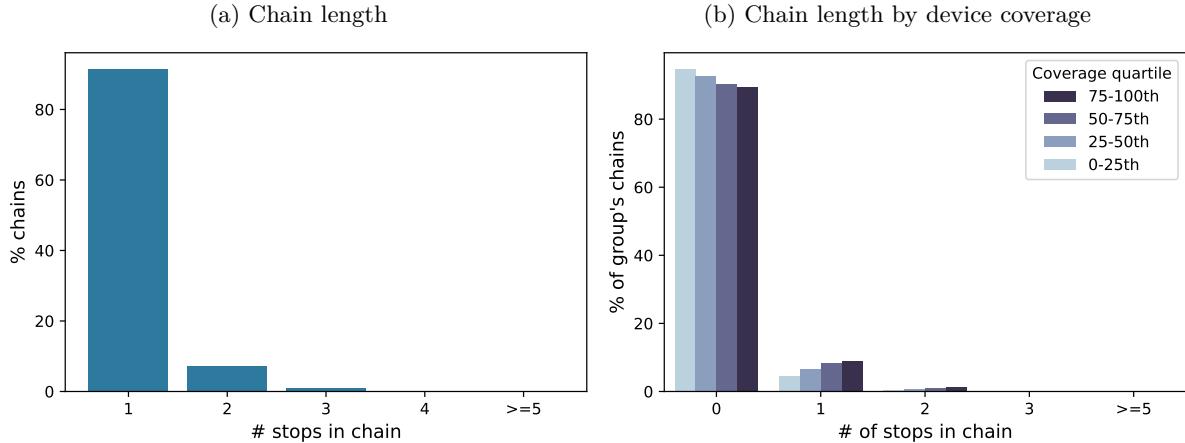
Figure A.1 Panel (a) plots the distribution of number of POIs visits in a trip – 91.5% of trips consist of a single stop. Panel (b) plots the distribution again, this time breaking out by quartiles device coverage. I define coverage at the device-week level as the sum of all time observed in the given week – on average, the top quartile of device-weeks are observed for about 83% of all minutes of the week, while the bottom quartile are observed for about 11%. For devices in the top quartile of coverage, 89.5% of chains consist of a single stop, compared to 94.7% of trips in the bottom quartile.

Table A1: Categories and subcategories of amenities

Sub-category name	NAICS codes	Top 100 CBSAs		Driving time (home)		
		# POIs	# choices	p25	p50	p75
<b>Restaurants</b>						
Full-service	722511	315056	13403993	10.1	15.7	23.8
Limited-service	722513	146561	4004882	9.4	15.0	23.4
Cafes; snacks	722515	95190	2547540	9.2	14.9	23.3
Drinking places	722410	47643	2261690	10.0	16.7	25.8
<b>Shops</b>						
Gas stations, convenience	4471, 445120	99360	4676792	9.0	15.3	24.9
Clothing, shoes	4481, 4482	95253	2199797	10.4	16.5	25.1
Grocery stores	445110, 4452	81533	4146230	8.0	12.4	19.5
Sporting goods, hobby, music	4511	54032	1236785	10.5	16.3	24.6
Automobile dealers	4411, 4412	50172	2332453	12.0	18.6	27.4
Furniture, furnishings	4421, 4422	48845	1197002	11.2	17.5	26.4
Building materials, gardening	4441, 4442	41531	2746107	9.8	14.3	21.1
Beauty, glasses, personal care	446120, 446130, 446190	40071	1635234	9.9	14.8	21.9
Pharmacies	446110	39857	1849610	8.4	13.6	21.2
Beer, wine, liquor, tobacco	4453, 453991	38261	730683	8.9	14.4	22.3
Electronics, appliances	4431	29603	603816	10.3	16.4	25.0
General merch., warehouse clubs	4523	26072	6032372	9.4	13.8	19.8
Jewelry, luggage, leather	4483	25605	646372	9.6	15.7	24.3
Malls	N/A	25145	78953287	9.7	14.5	21.6
Used merchandise	4533	23245	630241	9.4	14.8	22.5
Office supplies	4532	22407	440852	10.4	16.7	26.0
Department	4522	10324	1235584	9.4	13.6	19.8
Books, news	4512	8335	207113	9.9	15.7	23.8
<b>Personal services</b>						
Salons, barbers	8121	287677	6271867	9.0	14.1	21.4
Mechanics	811111, 811112, 811113, 811118, 811121, 811122, 811191, 811198	142078	1011801	9.0	15.1	23.9
Religious organizations	8131	127007	3985066	9.4	15.0	22.2
Other health practitioners	6213	123280	2056482	10.2	15.9	23.7
Banks	5221, 5223, 5231	119383	2518547	8.6	13.9	22.1
Dentists	6212	112700	1742965	9.8	15.4	23.2
Telecom	5173	31577	445421	9.5	14.7	22.2
Accounting, taxes, payroll	5412	30317	380276	9.9	15.7	23.9
Hospitals	6221, 6214	29705	10040470	11.9	18.3	27.1
Car washes	811192	15835	191638	8.7	13.8	22.3
Drycleaning, laundry	8123	11258	211991	7.3	12.3	19.8
<b>Entertainment</b>						
Fitness centers	713940	78147	4839492	9.1	14.3	21.9
Parks	712190	76062	16955117	7.9	14.7	25.1
Other amusement, recreation	713990	16312	679272	9.3	16.1	25.5
Golf courses	713910	9847	6833737	5.5	14.0	24.6
Museums, zoos, gardens	712110, 712130	7526	252688	13.0	21.2	32.0
Libraries	519120	7220	480875	7.9	12.4	19.3
Movie theaters	512131	6307	998776	11.2	16.4	23.6
Amusement parks	7131	4452	891530	14.1	23.3	36.6
Performing arts	7111, 7113	3411	847696	14.0	22.2	33.0
Spectator sports	7112	2999	1300992	15.1	23.4	34.1
Bowling	713950	2620	538098	10.9	16.2	23.5
Gambling	7132	1059	983991	15.4	23.4	35.5

*Note:* This documents the categories and subcategories of amenities, the NAICS codes used to identify each subcategory, and the total number within the top 100 CBSAs. NAICS codes are provided by SafeGraph. To identify malls, I use SafeGraph's category 'Lessors of Real Estate.' Driving time is between the home block group of the device observed visiting the establishment and the block group of the establishment.

Figure A.1: Distribution of # stops in chain



*Note:* These figure document the distribution of the number of stops at different POIs within a single trip, where a new trip begins if a device stops at home or work or it has been at least 2 hours since the device left the previous POI.

## A.2 Estimating individual income

I use parcel-level data to estimate whether or not a device residing at that parcel is above the median income for the CBSA ('higher income').

At a high level, the procedure is as follows:

1. Match each device to a home parcel using data on parcel polygons and Corelogic parcel characteristics
2. Use data on historical transactions to estimate the 2019 market value of all houses
3. Model the probability that a device-holder is above median income based on their home's market value, location, and characteristics using the 2019 5-year ACS microdata
4. Bayesian-update the probability of higher income using the distribution of home block group income in the 2019 5-year ACS tabulations

Each step is described in more detail below.

### A.2.1 Matching devices to houses

I match each device to residential parcel polygons provided by LandGrid by spatially joining each set of home coordinates to parcel data. I use a strict match and only include those home coordinates that fall within the bounds of a residential parcel.

Next, I match LandGrid parcels to housing characteristics provided by Corelogic. There are two main data sources from Corelogic: characteristics of a given house from tax assessments filed with local governments ("Tax") and details of each sale of a house from the deed filed during sale ("Deed"). Most of the Tax data is from the 2018 tax year, while the Deeds extend through June

2019. To match Landgrid parcels to Corelogic details, I use the coordinates provided by Corelogic and again spatially join to the polygon from Landgrid.

In total, 73.2% of device homes match to a parcel with Corelogic characteristics. The majority of the remaining devices were not successfully matched to any residential parcel (e.g., their ‘home’ coordinates may have been on the street if home identification was noisy for that device).

### A.2.2 Estimating market value of homes

To predict the value of a given device’s home, I build a model for estimating the 2019 market value of all housing units in the US (not just those with a matched GPS device). The model is trained on data on actual housing sales from Corelogic, then used to predict out-of-sample using the location and characteristics of all homes, even those for which I never observe a sale.

I focus on predicting the value of single-family houses, townhouses, and condos; large rental apartment buildings are, for now, set aside. There are 104.1 million properties in the Tax data for which I will ultimately predict the 2019 market value.

To build a training sample, I match each house in the Tax data to any sales between 2010–2019 in the Deed data. I restrict the sales data to arms-length transactions and remove all homes purchased by an owner with ‘LLC’ in the name, as they often involve purchases of many units that have sale prices corresponding to the total purchase rather than each unit. Finally, I further restrict to houses sold at prices between \$10,000 and \$25,000,000. The final data include 27.99 million sales.

The Tax data include many characteristics of the property, including: number of bedrooms and bathrooms; living, garage, basement, and land square footage; the style of the building, when it was built, and whether refurbished; a measure of quality (e.g. ‘Fair’); what type of view the property has (e.g., ‘Mountains’) and other ‘location influences’; whether the property has a pool; the type of air conditioning, heating, fuel, framing, walls, sewage, water, roof, and floor; and the precise location. Some of these variables are frequently missing – in all such cases, I replace the missing values with an indicator for missing data. I also discretize continuous variables such as living square footage into deciles.

To predict the market value of all homes, I train a deep neural net (DNN) on 80% of all sales, setting aside the remaining 20% to evaluate model performance. First, I turn all of the characteristics of the property and the time of the sale into features. Many of these characteristics are high dimensional; for example, there are 87 unique styles of houses and 287 types of exterior walls. For sale year, sale month, property type, decade built, decile of living square feet, and whether there is a pool, I encode the values using dummy variables. For the remaining characteristics described above, I use embeddings. Embeddings are an alternative to dummy variables – rather than mapping a unique value to a vector of 0s and a single 1, they map each unique value to a vector of continuous values whose weights are learned through the model training. Using embeddings

helps reduce the dimensionality and allow the model to discover similarities between categories (e.g., similar types of building styles). As a rule of thumb, I define the embeddings to have dimensionality of half the number of unique values of a column, with a max of 500. I also use embeddings for the Census tract in which the house is located. By using embeddings, the model can ‘discover’ that two tracts (perhaps sharing a border) are quite similar; if the tracts are similar in their relationship with sale price, they will be close in the embedding space as well (i.e. their embedding vectors will have similar values).

Figure A.2 describes the model architecture. I use 5 hidden linear layers with rectified linear units (ReLUs) as activation functions between each layer and batch normalization to increase stability. I apply a sigmoid function to the final output to restrict estimates to be between 70% of the lowest observed sale price and 130% of the highest observed sale price. The model is built using PyTorch and run on GPUs. I train the model on each state’s data separately for 100 epochs with batches of 256 sales (i.e. 100 runs of the full data through all steps of the model architecture by batch, with parameters updating after each batch). I use an Adam optimizer with a cosine annealing learning rate, such that the learning rates starts high and then decreases following a half-cosine curve.

Figure A.2: Model architecture

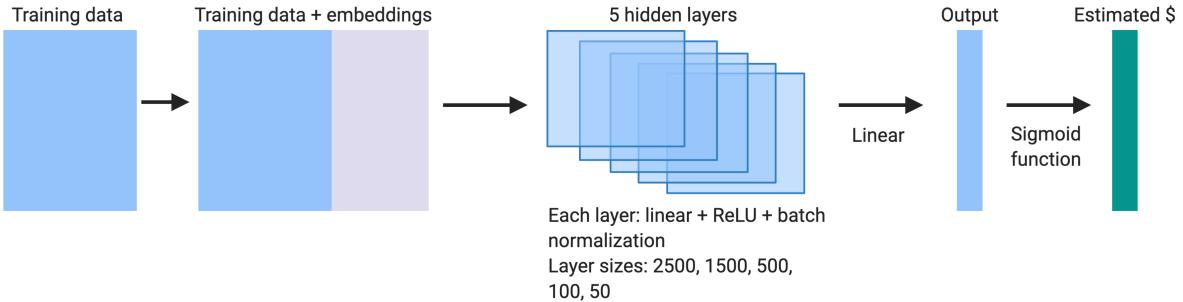


Table A2 compares the performance of the DNN on the 20% validation sample compared to the performance of a simple linear regression model for sales in California. For the linear regression, features are encoded as dummy variables instead of embeddings and, for computational purposes, zipcodes are used rather than the more granular Census tracts. The first two columns display the root mean squared log error in training and testing samples. The similarity between training and testing samples suggests the model is not overfit, and the lower value compared to the linear regression indicates the superior fit of the DNN with embeddings over the linear regression with all dummy columns. As a second measure of model fit, I look at the percent of homes for which the predicted sale amount is in the same quartile, decile, or 25-tile of the actual sale amount. For 79.20% of homes, the predicted value is in the same quartile (compare to 51.64% for the linear regression model and 25% if randomly assigned). Figure A.3 provides a visualization of how well the model predicts the correct decile; rows are decile of predicted sale amount, columns are deciles

of actual sale amount, and the color corresponds to the number of sales. A model that could perfectly separate houses into deciles would be a diagonal line of dark squares. In general, the DNN is far closer to this ideal than the linear regression and is only infrequently off by more than a decile.

Table A3 reports similar statistics for all other states modeled using the DNN. On average, 65.38% of houses in the testing sample have predicted sales amounts that are in the same within-state quartile as the true sale amount.

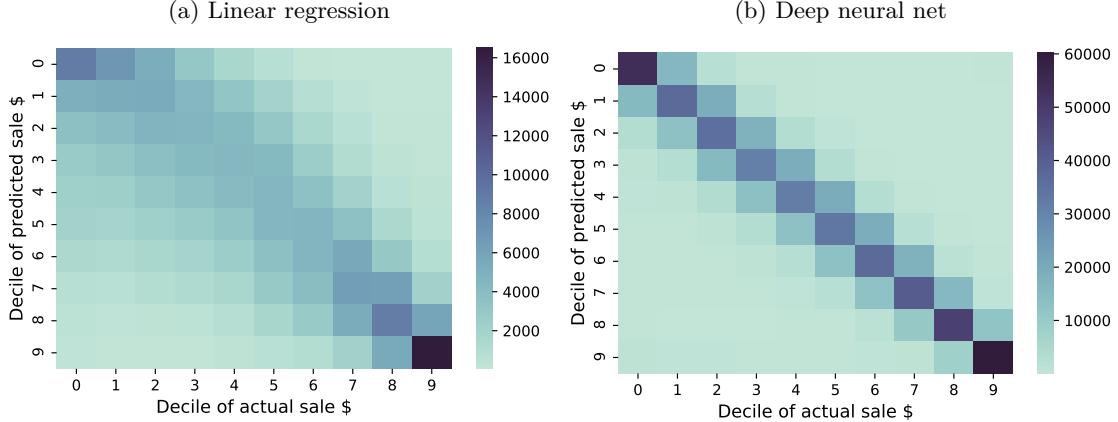
Finally, for all 104.1 million properties in the Tax data—even those that did sell—I predict the market value based on the property characteristics and assuming a June 2019 sale date.

Table A2: Comparing models: California

	Root mean squared log error		% in correct p-tile		
	Training	Testing	Quartiles	Deciles	25-tile
Deep neural net	0.6055	0.6161	79.20	54.32	28.72
Linear regression	0.7145	0.7130	51.64	25.32	11.36

*Note:* This table compares model performance between the neural net with embeddings and a linear regression with all categorical variables represented as dummies. The results shown are just for sales in California. ‘Training’ refers to the 80% of data used to train the model, ‘testing’ refers to the 20% set aside. The % in correct p-tile is based on the testing sample and is defined as the percent of houses for which the p-tile of predicted sale amount is the same as the p-tile of actual sale amount.

Figure A.3: Deciles of true v. estimated sale amounts: California



*Note:* These figure shows the out-of-sample quality of different models for California by highlighting the percent of predicted sale value deciles that match the actual sale value.

### A.2.3 Estimating whether above median income using home characteristics

I use data from the 2019 5-year ACS Public Use Microdata Sample (PUMS) to model the relationship between housing characteristics and income. I then use this model to predict whether a device

Table A3: DNN performance: all states in sample

State	Sample size		Root mean squared log error		% in correct p-tile		
	Num. homes	Num. sales	Training data	Testing data	Quartiles	Deciles	25-tiles
Alabama	1,888,962	387,858	0.6859	0.7070	51.39	26.11	12.10
Alaska	236,172	55,985	0.4057	0.4393	62.97	35.57	16.96
Arizona	2,236,852	891,688	0.4986	0.5006	67.90	39.79	19.77
Arkansas	1,194,026	273,041	0.7086	0.7141	56.13	29.69	13.69
California	9,568,480	3,033,246	0.6055	0.6161	79.20	54.32	28.72
Colorado	1,887,238	707,054	0.4843	0.5008	74.09	47.23	23.89
Connecticut	1,119,230	281,333	0.5907	0.6044	67.38	39.99	20.07
Delaware	344,763	81,055	0.5627	0.5539	61.95	33.59	15.97
District of Columbia	153,194	48,502	0.4193	0.4122	70.74	42.34	21.03
Florida	7,427,890	2,729,955	0.7170	0.7273	62.61	34.39	16.88
Georgia	3,588,908	927,570	0.8304	0.8452	53.31	27.73	13.29
Hawaii	435,009	116,236	0.4579	0.4869	73.49	46.99	23.63
Idaho	649,577	158,018	0.5303	0.5226	54.60	27.38	11.96
Illinois	4,016,897	1,140,526	0.6637	0.6896	64.39	35.82	17.33
Indiana	2,314,467	555,692	0.7905	0.7926	49.44	25.00	11.79
Iowa	1,251,933	325,829	0.8950	0.8890	51.79	25.28	11.29
Kansas	949,340	170,364	0.5293	0.5331	66.54	37.79	17.58
Kentucky	1,706,493	375,699	0.6076	0.6183	58.41	31.07	14.31
Louisiana	1,650,607	319,622	0.7391	0.7551	51.40	26.18	12.03
Maine	539,490	77,963	0.5697	0.5795	51.78	25.87	11.47
Maryland	1,919,137	544,304	0.5175	0.5165	70.34	42.88	21.89
Massachusetts	2,011,990	575,758	0.4893	0.4904	70.40	42.35	20.79
Michigan	3,925,464	996,063	0.7569	0.7735	55.14	28.70	13.32
Minnesota	1,917,659	593,594	0.5973	0.6096	59.56	32.23	15.15
Mississippi	1,205,515	72,095	0.6694	0.6533	58.14	28.76	12.25
Missouri	692,697	88,942	0.5524	0.5698	56.43	30.63	14.27
Montana	386,693	77,779	0.5125	0.5168	54.26	27.72	12.38
Nebraska	661,336	196,668	0.5055	0.5245	66.86	39.60	19.29
Nevada	919,364	387,895	0.4890	0.4842	73.63	46.36	23.46
New Hampshire	527,069	132,949	0.5176	0.5277	58.82	31.60	14.97
New Jersey	2,548,556	730,375	0.5079	0.5276	68.90	40.05	19.53
New Mexico	696,968	113,863	0.6867	0.6849	48.15	23.34	10.53
New York	4,594,929	1,193,795	0.6376	0.6486	71.73	42.34	19.92
North Carolina	3,922,110	923,345	0.6156	0.6287	65.55	37.58	18.18
North Dakota	219,117	61,043	0.6745	0.6766	53.03	24.86	10.53
Ohio	4,045,159	1,090,874	0.5811	0.5892	61.93	34.37	16.57
Oklahoma	1,562,997	368,883	0.6125	0.6380	63.22	35.63	17.05
Oregon	1,355,705	452,541	0.4698	0.4800	69.31	41.16	19.73
Pennsylvania	4,402,072	1,064,641	0.5576	0.5705	62.96	35.42	17.40
Rhode Island	329,573	84,482	0.5070	0.5126	61.85	34.13	16.40
South Carolina	2,008,807	493,926	0.7231	0.7240	52.96	27.05	12.85
South Dakota	272,046	34,026	0.6346	0.6490	56.81	28.38	12.87
Tennessee	2,702,028	735,374	0.6359	0.6534	64.02	36.33	17.84
Texas	8,281,129	1,727,422	0.8504	0.8468	60.17	32.78	15.03
Utah	862,770	257,946	0.4857	0.4937	61.81	33.30	15.45
Vermont	254,667	76,678	0.8029	0.8178	48.19	23.49	10.52
Virginia	3,154,047	803,333	0.6482	0.6639	75.15	49.39	26.00
Washington	2,370,843	772,662	0.5926	0.6002	70.64	42.72	21.29
West Virginia	726,247	41,956	0.9613	0.9440	51.59	25.03	9.33
Wisconsin	2,209,251	605,939	0.7468	0.7511	50.94	25.14	11.17
Wyoming	209,870	35,451	0.5115	0.5225	57.17	30.45	14.37
All	104,055,343	27,991,838	0.6422	0.6522	64.62	37.50	18.39

*Note:* This table documents performance for each US state for which I have housing data. ‘Training’ refers to the 80% of data used to train the model, ‘testing’ refers to the 20% set aside. The % in correct p-tile is based on the testing sample and is defined as the percent of houses for which the p-tile of predicted sale amount is the same as the p-tile of the actual sale amount.

is above median income for their home CBSA using the characteristics of the parcel in which they reside as well as the overall income distribution of their block group.

I first harmonize the ACS and Corelogic characteristics, such that a model can be trained on the ACS data and evaluated on Corelogic data. I use up to four characteristics of each home: its location (Public Use Microdata Area), the number of units in the building, the decade built, and its estimated home value. In the ACS, households self-report their home value ('valueh') or the rent they pay ('rentgrs'). In Corelogic, I observe only estimated home value and not whether the parcel is occupied by the landlord or renters. To address this, I use the within-CBSA decile of either home value or gross rent (whichever is available) instead of the actual estimated market value. This approach is equivalent to simply using the decile of home value if rents are set as a CBSA-wide multiplier on the home value.

I use a common machine learning classifier, XGBoost, to estimate whether a household is above median income in the ACS data ([Chen and Guestrin, 2016](#)). I train three versions of this model, for varying levels of data availability: 1) using all four characteristics described above, 2) using just home value/rent deciles and the PUMA, and 3) using just the decade built, units, and PUMA. Each feature enters as a set of dummy variables, and I allow XGBoost to determine which interactions are important. I weight each household according to the household weights provided in the ACS. Evaluated on a holdout sample the three versions of the model correctly predict whether a household is above median income for 68.4%, 66.7%, and 65.1% of households respectively. The relatively low precision indicates that there remain many unobservables that contribute to household wealth beyond a home PUMA and housing characteristics, which motivates the value of using block group level income distributions to Bayesian-update the baseline predictions. I discuss this further below.

I take the model estimated on ACS data to the Corelogic data and predict the probability that residents of each parcel are above median income based on the available parcel characteristics. When the parcel does not have an estimated home value—e.g., for a large apartment building with many units—the probability that residents are above median income depends on the age of the building, number of units, and PUMA.

The ACS PUMS data includes only the PUMA of the household, but in the parcels data I can observe the exact location. To incorporate this information into the estimates of whether a device is higher income, I use the 2019 ACS 5-year block group income distributions to identify where above/below median income households live within the PUMA. I then use this to update the baseline probability that a parcel's residents are above median income using Bayes rule. In particular, for a tenant living in a parcel with characteristics  $x$  in block group  $g$ , I evaluate the probability they are above median income ( $H$ ) as

$$P[H | g, x] = \frac{P[H | x]P[g | H]}{P[H | x]P[g | H] + P[L | x]P[g | L]}$$

where  $L$  denotes below median income and  $P[g \mid H]$  is the within-PUMA probability that a household of type  $H$  lives in block group  $g$ .

Finally, I use the mapping of parcels to devices to assign a probability of being above median income to each device. ‘Higher income’ devices are those for which the probability of being above median income, based on their home parcel, is at least 50%. In total, I am able to generate a prediction of device type for 97.8% of the devices matched to a Corelogic parcel.

#### A.2.4 Estimating a continuous measure of income

For robustness tests, I define higher and lower income devices based on being in the top or bottom tercile of the within-CBSA income distribution. To divide devices into terciles, I estimate their income in levels following a similar procedure to the above. First, I estimate the relationship between household income ( $y$ ) and housing characteristics ( $\mathbf{x}$ ) in the ACS. I use the same features as before but now use Ordinary Least Squares instead of an XGBoost classifier for estimation. Using the estimated model, I compute the predicted income as well as the variance of the prediction error for each parcel in Corelogic. Assuming the error is normally distributed, this provides an initial probability density function  $p(y \mid \mathbf{x}_i)$  for parcel characteristics  $\mathbf{x}_i$ .

Next, I update each income estimate using the distribution of household income for the corresponding block group. This will push the initial income estimates towards the average income of the corresponding block group, with larger changes for more noisy estimates of income. The ACS provides block group counts of households in 16 different bins of household income, ranging from \$0-\$10,000 to >\$250,000. Define  $h(b \mid \mathbf{x}_i, w_i)$  as the probability a household’s income is in bin  $b \in B$  given parcel characteristics and home block group  $w_i$ . Under the assumption that the block group income distribution  $g(b \mid w_i)$  and the estimated distribution based on parcel characteristics  $p(y \mid \mathbf{x}_i)$  are independent,  $h(b \mid \mathbf{x}_i, w_i)$  is given by

$$h(b \mid \mathbf{x}_i, w_i) = \frac{g(b \mid w_i) \int_b^{\bar{b}} p(y \mid \mathbf{x}_i) dy}{\sum_{b' \in B} \left[ g(b' \mid w_i) \int_{\underline{b}'}^{\bar{b}'} p(y \mid \mathbf{x}_i) dy \right]}$$

where  $\underline{b}$  and  $\bar{b}$  are the lower and upper household income bounds of bin  $b$ . The final income estimate is then a weighted average of bin midpoints with weights corresponding to the probability that the residents of a given parcel have household income within that bin:

$$\hat{y}_i = \sum_b \left[ \frac{(\bar{b} - \underline{b})}{2} h(b \mid \mathbf{x}_i, w_i) \right] \quad (\text{A.1})$$

Finally, to split devices into terciles, I define within-CBSA cutoffs using the distributed of estimated income.

### A.3 Evaluating sample quality

Evaluating the sample quality is complicated by the lack of any ‘ground truth’ about the demographics of devices in the data. To provide some sense of coverage, I compare the inferred home locations of device holders to the true population, using the 5-year 2019 American Community Survey (ACS).

First, I compare the spatial distribution of the device sample to that of the total population. While homes are re-assigned each quarter that a device is present in the data, for the purposes of comparing to the Census I use only unique homes for each device (generally a single home). Figure A.4 maps the ratio of devices to population for all the US counties in the sample. While there is better coverage in some areas (e.g., Florida) than others (e.g., New England), in general the sample includes 5-10% of a county’s residents, with an average coverage of 7.07%.

Next, I zoom in to the Census block group level—each of which covers approximately 1,500 residents—and examine whether devices are disproportionately sampled from block groups with certain demographics. To do so, I divide block groups into within-CBSA deciles by various characteristics weighting by the population of each block group. Then, I compute the number of devices with homes in block groups corresponding to each decile.

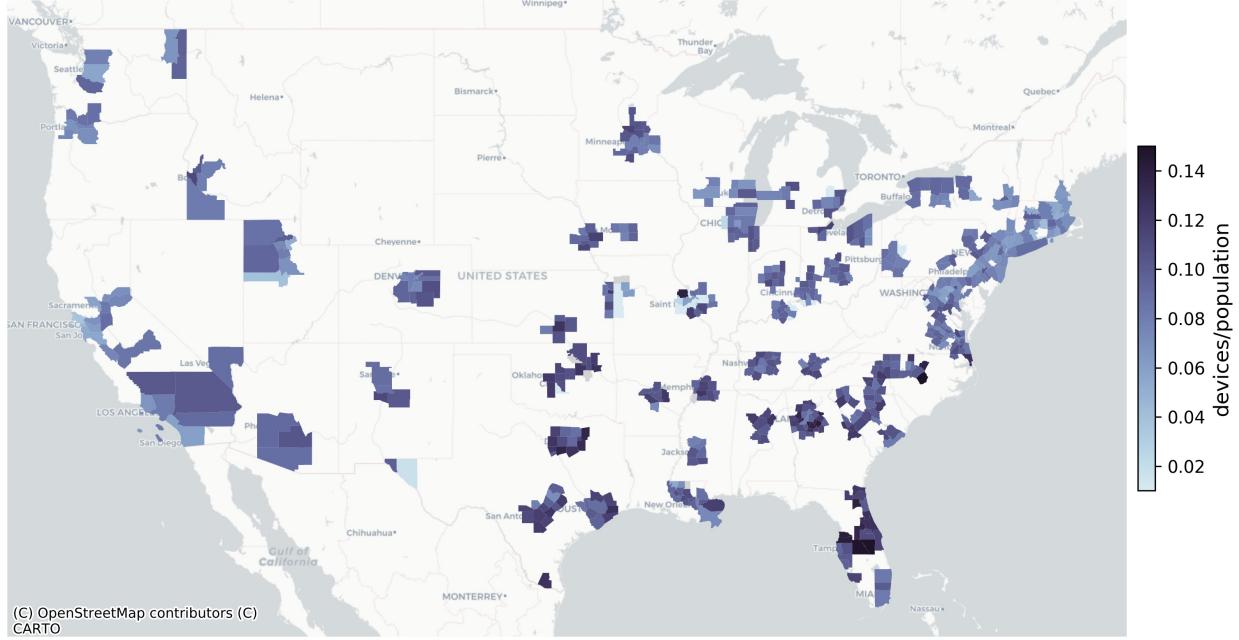
Figure A.5 plots the fraction of devices from each decile for block group population density, median household income, share of residents with a college degree, and the share of residents who identify as white. If the sample were perfectly uncorrelated with these characteristics, then 10% of devices would come from each decile (corresponding to the dashed horizontal line). Instead, devices in the sample, on average, come from lower income, less educated, and slightly less white block groups. One contributing factor is likely that Android devices are over-represented in the GPS data. While the data I have do not include the device type, anecdotally vendors have said that iPhones are only 35-50% of their sample despite making up about 60% of US smartphones. According to a survey, iPhone owners have incomes that are, on average, 43% higher than Android owners ([Slickdeals, 2018](#)). Compared to the sample of devices used in [Couture et al. \(2021\)](#), this sample is more evenly sampled across population densities but are less representative of neighborhood income, education, and fraction white. These discrepancies in coverage motivate the use of sample weights for all analyses.

## B Model appendix

### B.1 Overview of tensor-based estimation frameworks

PyTorch and its primary competing framework, Tensorflow ([Abadi et al., 2016](#)), are most commonly used to estimate complex neural networks, but can be applied to the estimation of many economic models. These frameworks offer a number of advantages for when the data are too large to fit in memory and the optimization process is too costly for more standard approaches, most notably

Figure A.4: Device homes as fraction of ACS population



*Note:* This figure maps the number of devices as a fraction of the total population for each county within a CBSA in the sample.

methods for estimating the model while holding only small batches of data in memory at any time as well as easy portability to estimating on Graphical Processing Units (GPUs).

The use of PyTorch and Tensorflow to estimate large economic models is relatively new to the economics literature.<sup>24</sup> At their core, these frameworks are solving a standard numerical optimization problem – given some objective function, find the set of parameters that maximizes (or minimizes) its value. In my setting, the objective function is the log-likelihood of observing visits to specific establishments, but similar techniques can be applied to a GMM criterion.

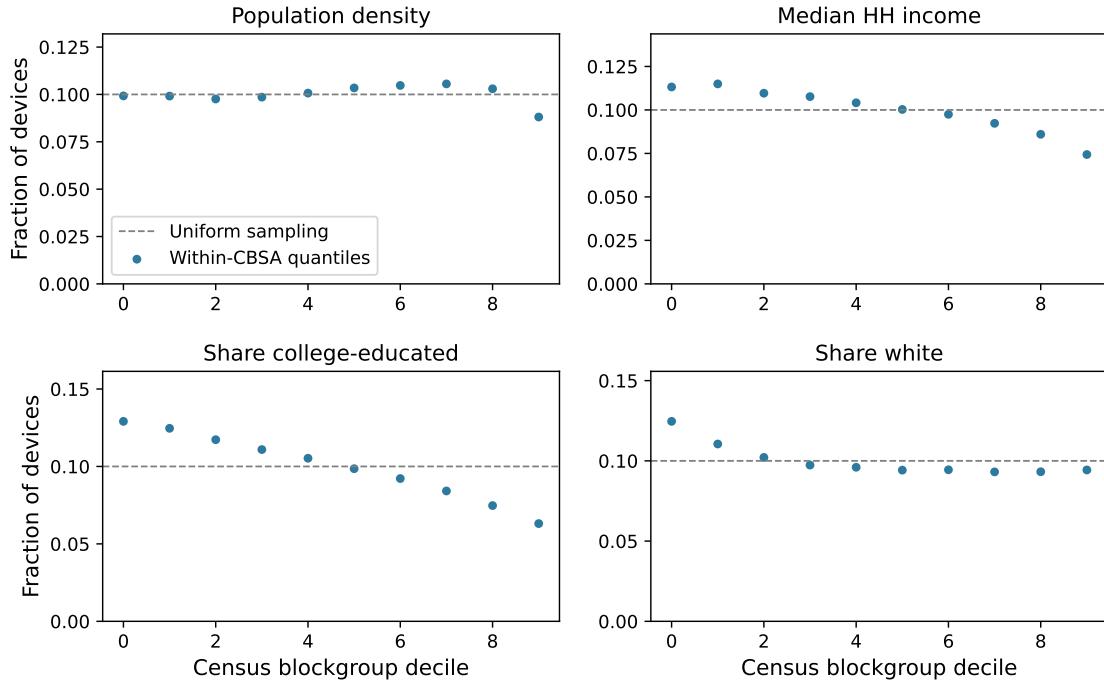
PyTorch and Tensorflow offer a number of advantages over standard approaches to optimization. First, these frameworks can easily scale to larger-than-memory data through ‘mini-batch gradient descent’ learning algorithms. ‘Mini-batch’ refers to reading in small subsets of data into memory at a time, evaluating the objective function on this mini-batch, updating the parameters, then proceeding to the next mini-batch.

Second, for updating the parameters these frameworks include many optimizers that rely on only the gradient and avoid ever computing the Hessian. For models where the parameter space is large, computing the Hessian becomes computationally infeasible. There are many gradient-based optimizers designed to squeeze as much as possible out of the gradient alone; in my case, I use

---

<sup>24</sup>The only example I have encountered in the economics literature is Lewis et al. (2021), which discusses an MLE estimator for BLP built in PyTorch.

Figure A.5: Device sampling by block group characteristics



*Note:* This figure shows the sampling of devices from block groups relative to a uniform sampling. Dots above the line indicate that more than 10% of devices come from that block group. Deciles are calculated either within-CBSA.

Adam, a popular optimizer that uses adaptive learning rates for each parameter and a weighted average of gradients from past mini-batches to avoid over-shoot (Kingma and Ba, 2017). The mini-batch gradient descent learning algorithms make it possible to estimate models when both the data are too large to store in memory and the parameter space is too large to efficiently compute the Hessian.

Third, these frameworks use ‘auto-differentiation’ to compute the gradient without being fed the analytical form of the gradient or relying on numerical approximations. This is particularly useful when training large neural networks, where the gradients are difficult to compute analytically and numerical approximations accumulate substantial error; however, even in the case of a more simple model, auto-differentiation can save the researcher time. For all its benefits, auto-differentiation is also a limitation of these frameworks, as they can only optimize objective functions that are constructed out of auto-differentiable building blocks. These building blocks include all the standard matrix operations and elementary functions that can be linked together by the chain rule. Certain problems, such as those that involve numerical integration or for-loops, may not be as amenable to auto-differentiation; however, these frameworks are still young and the space of models they can estimate is still being explored.

Finally, estimation can be easily run on either CPUs or GPUs. GPUs, which are traditionally

used for powering graphics, excel at executing many small tasks (such as matrix multiplications) in parallel. They offer substantial speed improvements over CPUs.

## B.2 Estimation of lower level of model

For each subcategory  $B_m$ , I use data on all visits to establishments within this subcategory to estimate the establishment-level taste parameters and the disutilities to driving time from home and work.<sup>25</sup>

I use a fixed choice set consisting of all establishments in a subcategory for which I observe at least 5 visits in the data. For many subcategories, the choice sets are large; for example, there are over 12,000 full-service restaurants in Los Angeles. I estimate the model by maximizing the log-likelihood, based on the conditional probability of consuming at a specific establishment within subcategory  $B_m$  defined in Equation 4.4. For estimation, I use a mini-batch gradient descent algorithm built in PyTorch and run it on a virtual machine with 16 CPUs, 104GB of memory, and 1 Nvidia T4 GPU.

The estimations steps are outlined in Algorithm 1. Much like for the model to predict house prices, I use an Adam optimizer with an initial learning rate of 5e-3, a cosine annealing learning rate scheduler<sup>26</sup>, and a convergence tolerance of 5e-6. Adam uses only the current and historical gradients to update parameters and adjusts the learning rates for each parameter separately (hence the  $\gamma$  and  $\kappa$  superscripts on  $\eta$  in Algorithm 1). I weight the loss function by a given device’s sample weight. To represent the fixed effects for the  $C$  different establishments in the choice set, I use a  $1 \times C$  vector of embeddings rather than a one-hot encoding (‘dummy variables’). While embeddings are usually used to model some latent space—see, for example, word2vec (Mikolov et al., 2013)—they can also be used as a sparse way to represent high dimensional fixed effects. Given the size of the choice sets, standard one-hot encoding would be infeasible in my setting.

There are a few cases—especially for the ‘Gambling’ and ‘Museums, zoos, gardens’ subcategories—where a CBSA will have zero establishments and is therefore excluded from estimation.

## B.3 Estimation of upper levels of model

I estimate the upper levels of the model using the aggregate number of visits a given device made to each subcategory  $B_m$  at each time of week  $t$  within a quarter. The likelihood contribution of a

---

<sup>25</sup>There are a few cases of driving times being unreasonable long for within-CBSA trips, likely due to errors in the router, so I top-code driving times at 90 minutes and include an indicator for whether a driving time was over 90 minutes.

<sup>26</sup>For more on the Adam optimizer, see Kingma and Ba (2017). Loshchilov and Hutter (2019) show that, even with adaptive learning rate optimizers like Adam, it can still improve estimation to use a scheduled learning rate multiplier, such as cosine annealing.

---

**Algorithm 1** Lowest-level estimation: mini-batch gradient descent

---

- 1: Initialize  $1 \times C$  embedding vector to store establishment-level taste parameters ( $\vec{\gamma}$ ) and a  $1 \times 12$  vector of driving time from work and home disutility parameters ( $\vec{\kappa}$ )
- 2: **while** not converged and  $epoch \leq maxEpochs$  **do**
- 3:   Divide sample into mini-batches, each with 128 observed choices
- 4:   **for**  $b \in$  mini-batches **do**
- 5:     Expand each choice into the full choice set; compute driving time to each establishment
- 6:     Compute predicted choice,  $\hat{y}$ , under current parameters (forward propagation)
- 7:     Compute weighted cross-entropy loss for  $\hat{y}$  and true choices  $y$ ,  $\mathcal{L}(\hat{y}, y, deviceWeights)$
- 8:     Compute gradients using backpropagation ( $\nabla \vec{\gamma}, \nabla \vec{\kappa}$ )
- 9:     Update parameters using learning rates  $\eta_t^\gamma$  and  $\eta_t^\kappa$

$$\begin{aligned}\vec{\gamma}_{t+1} &\leftarrow \vec{\gamma}_t - \eta_t^\gamma \nabla \vec{\gamma}_t \\ \vec{\kappa}_{t+1} &\leftarrow \vec{\kappa}_t - \eta_t^\kappa \nabla \vec{\kappa}_t\end{aligned}$$

where  $t$  indexes ‘time’ as the current (epoch, mini-batch) pair

- 10:   Check for convergence of loss function
  - 11:   Shuffle the data
- 

vector  $\vec{n}_{iqt}$  for arrival rate  $\lambda_{kt}$  is given by

$$\begin{aligned}P[\vec{n}_{iqt}] &= \sum_{a=0}^{\infty} P[a \text{ arrivals} | \lambda_{k(i)t}] * P[\vec{n}_{iqt} | a] \\ &= \sum_{a=0}^{\infty} \left[ \left( \frac{\lambda_{k(i)t}^a \exp(-\lambda_{k(i)t})}{a!} \right) * \left( \frac{a!}{n_{iqt0}! \prod_{B_m} n_{iqtB_m}!} \right) * \left( \prod_{B_m} P_{itB_m}^{n_{iqtB_m}} \right) \left( P_{it0}^{n_{itq0}} \right) \right] \\ &= \sum_{a=0}^{\infty} \left[ \left( \frac{\lambda_{k(i)t}^a \exp(-\lambda_{k(i)t})}{a!} \right) * \left( \frac{(a - n_{iqt0})!}{\prod_{B_m} n_{iqtB_m}!} \right) * \left( \prod_{B_m} P_{itB_m}^{n_{iqtB_m}} \right) \left( P_{it0}^{n_{itq0}} \right) \right] \quad (\text{B.1})\end{aligned}$$

where  $n_{iqtB_m}$  indexes the  $B_m$  subcategory of  $\vec{n}_{iqt}$ ,  $P_{itB_m}$  is computed according to Equation 4.7,  $P_{it0}$  is the implied probability of consuming the outside option, and  $n_{itq0}$  is the implied number of outside choices for a given number of arrivals  $a$ .<sup>27</sup> To approximate the infinite sum, I compute the summation for up to 100 arrivals – higher numbers of arrivals have probabilities that are extremely small (i.e. too small for even a 64-bit floating point value to store). As only relative differences are identified, I normalize the first coefficient for each vector of intercepts to be zero.

Devices are often observed in the data for a subset of the total quarter – without accounting for this, the arrival rate would be conflated with device coverage. To address this issue, I compute the number of unique days a device is observed in each quarter—whether or not they go to a POI that day—and substitute in an individual level arrival rate  $\lambda_{iqt}$  in Equation 4.10, with  $\lambda_{iqt} = \lambda_{k(i)t} * numDays_{iqt}$  where  $numDays_{iqt}$  is the number of days device  $i$  is observed in quarter  $q$  at

---

<sup>27</sup>These are given by  $P_{it0} = 1 - \sum_{B_m} P_{itB_m}$  and  $n_{itq0} = a - \sum_{B_m} n_{itqB_m}$

time of week  $t$ . As such,  $\lambda_{igt}$  can be interpreted as a daily arrival rate for each time of week (e.g., number of arrivals on a weekday evening).

For the upper levels of the model, I again use a mini-batch gradient descent learning algorithm built with PyTorch to estimate the model.<sup>28</sup> I use batch sizes of 128 device-quarters, an Adam optimizer, and an initial learning rate of 1e-3. To help avoid finding only local maxima—the objective function of a nested logit is not globally convex—I use a cosine annealing learning rate scheduler that increases the learning rate every so often to try to jump out of any local maxima. In testing, different initializations of the parameters led to similar final estimates. I train the model for a maximum of 10 epochs, with a check for convergence. I use fewer epochs than in the lower level because these models tend to converge far more quickly—usually after only a few epochs—thanks to the smaller parameter space.

## C Additional results

### C.1 Preferences, exits, and gentrification

In this section, I show how establishment-level preferences affect the probability that an establishment exits a neighborhood, split by whether or not the neighborhood is ‘gentrifying.’ I combine data from SafeGraph and Yelp to infer whether an establishment closed. Yelp data is sparsely available for many categories of amenities, so I focus specifically on restaurants, where the vast majority of establishments have a Yelp page. SafeGraph flags that an establishment has likely closed if it disappears from their underlying data sources for multiple months in a row. This method is imperfect, so I additionally check whether the establishment is still open in Yelp.<sup>29</sup> To be conservative, I exclude restaurants that either SafeGraph says are closed but Yelp says are open or that SafeGraph says are open but that do not match to a Yelp page. I label as ‘closed’ only those restaurants that SafeGraph labels as closed and Yelp does not list as still open. The final sample include 274,486 total restaurants, of which 30,175 closed (11%).

Next, I label each neighborhood as gentrifying or not. Gentrification is an ambiguous and often politicized term, originally coined by Ruth Glass to describe the entry of educated, middle-class residents (the ‘gentry’) to poor, working quarters of London (Glass, 1964). I follow Glaeser et al. (2020) and use a discrete definition of gentrification based on a zip code’s initial condition and subsequent increases in either rent or education. I categorize each zip code in a CBSA as one that could *potentially* gentrify if it was in the top 25th percentile of poverty in 2013, and, among these

<sup>28</sup>In this case, more traditional optimizers would also work well, because both the parameter space and data are far smaller.

<sup>29</sup>I pulled the Yelp characteristics used in analyses, such as price level, in October 2020. I checked whether an establishment is listed as open on Yelp in August 2021. I cannot rely exclusively on the Yelp data, as the Yelp API will generally only return details for establishments that are still open; when the API returns no result, it is unclear whether the establishment is closed or simply never had a Yelp page. For restaurants, 33% of the establishments that SafeGraph labels as having closed between January 2020 and August 2021 are still listed as open in Yelp.

zip codes, categorize a zip code as having actually gentrified if its increase in either median rent or fraction of residents with a college degree between the 2013 and 2019 5-year ACS was in the top half of these potentially gentrifying zip codes. For 89% of zip codes, these two measures of gentrification are equivalent.

I estimate the following linear probability model of exiting using OLS:

$$\text{Exit}_j = \text{gentrifying}_{z(j)} \times (\tilde{\gamma}_{jH} + \tilde{\gamma}_{jL}) + \epsilon_j \quad (\text{C.1})$$

where  $\text{Exit}_j$  is an indicator for whether establishment  $j$  exited,  $\text{gentrifying}_{z(j)}$  is an indicator for whether zip code  $z(j)$  is gentrifying and  $\tilde{\gamma}_{jH}$  and  $\tilde{\gamma}_{jL}$  are establishment-level DTEs.

The estimated coefficients are presented in Table C1. The first three columns use the definition of gentrifying based on increases in rent, while the latter three columns use the definition based on increases in education. The results are similar for the two measures.

Table C1 documents four facts about preferences, gentrification, and restaurant exits. First, restaurants that are liked by residents are less likely to exit. The effect is largest for higher income resident preferences – a 10 minute increase in the higher income DTEs of an establishment is associated with a 1.9 percentage point decrease in the probability of exit (columns 1 and 4). Second, the probability of exiting is larger for gentrifying neighborhoods, mirroring results in Glaeser et al. (2020). In a gentrifying neighborhood, the average establishment is 1.3-1.5 percentage points more likely to exit (columns 1 and 4). Third, resident preferences may matter slightly less in gentrifying neighborhoods than non-gentrifying neighborhoods, although the point estimates are small and generally insignificant (columns 2 and 5). Finally, establishments liked more by higher income residents than lower income residents are less likely to exit a gentrifying neighborhood. For an establishment where the gap between higher income and lower income residents is 10 DTEs, the point estimates suggest that it is 0.5-0.7 percentage points less likely to exit if it is in a gentrifying neighborhood than in a non-gentrifying neighborhood (columns 3 and 6).

The results provide some evidence of the endogenous response of a neighborhood's set of amenities to changes in the local demographics. When higher income residents value a restaurant more than lower income residents, it is less likely to exit a gentrifying neighborhood. However, these results should be interpreted with caution for a few reasons. First, the measure of exiting is imperfect and only available with any level of confidence for restaurants. Second, the January 2020 to August 2021 period overlaps almost entirely with the COVID-19 pandemic, which hit restaurants particularly hard. Finally, the measure of gentrification is based on historical changes to demographics; many establishments that were tailoring to lower income residents in these neighborhoods may have already exited by 2020, which would mute the magnitude of the coefficients.

Table C1: Preferences and exits: restaurants

	Rent gentrification			Education gentrification		
	(1)	(2)	(3)	(4)	(5)	(6)
Intercept	0.1040 (0.0007)	0.1040 (0.0007)	0.1047 (0.0007)	0.1038 (0.0007)	0.1037 (0.0007)	0.1044 (0.0008)
Gentrifying	0.0135 (0.0015)	0.0146 (0.0015)	0.0211 (0.0032)	0.0148 (0.0015)	0.0161 (0.0015)	0.0225 (0.0034)
$\widetilde{\gamma_{jH}}$	-0.0019 (0.0001)	-0.0020 (0.0001)		-0.0019 (0.0001)	-0.0020 (0.0001)	
$\widetilde{\gamma_{jL}}$	-0.0010 (0.0001)	-0.0011 (0.0001)		-0.0010 (0.0001)	-0.0011 (0.0001)	
Gentrifying $\times \widetilde{\gamma_{jH}}$		0.0002 (0.0002)			0.0001 (0.0002)	
Gentrifying $\times \widetilde{\gamma_{jL}}$		0.0003 (0.0002)			0.0004 (0.0002)	
$(\widetilde{\gamma_{jH}} - \widetilde{\gamma_{jL}})$			0.0002 (0.0001)			0.0002 (0.0001)
Gentrifying $\times (\widetilde{\gamma_{jH}} - \widetilde{\gamma_{jL}})$			-0.0005 (0.0002)			-0.0007 (0.0002)
$R^2$	0.0135	0.0135	0.0048	0.0136	0.0136	0.0049
N	277401	277401	277401	277401	277401	277401

*Note:* This table documents results from regressing whether a restaurant exited between Jan. 2020 and Aug. 2021 on the estimated preferences for that restaurant. I exclude restaurants that either SafeGraph says are closed but Yelp says are open or that SafeGraph says are open but that do not have a Yelp page. Gentrifying is defined at the zip code level based on being in the 25th percentile of poverty in 2013 and experiencing rent or education growth between 2013 and 2019 in the top half for these zip codes. All regressions include CBSA fixed effects. Standard errors are presented in parentheses.

## C.2 Tourist preferences

Many establishments cater to tourists and an increase to the number of tourists to an area of the city can affect its composition of amenities ([Almagro and Dominguez-Iino, 2022](#)). I re-estimate the model for visitors to cities, using their overnight location within the city as their ‘home.’ I find that tourist preferences for both establishments and neighborhoods are more correlated with higher income residents than lower income residents. The average within-CBSA, within-subcategory correlation in establishment level parameters is 0.68 for tourists and higher income residents and 0.61 for tourists and lower income residents. The corresponding correlations for NAQI values are 0.84 and 0.82. These results suggest that an increase in the number of tourists to a neighborhood is more likely to benefit higher income residents who have more similar preferences to the tourists, although the difference in correlations is small.

## C.3 Additional tables and figures

Figure C.1: Driving time coefficients

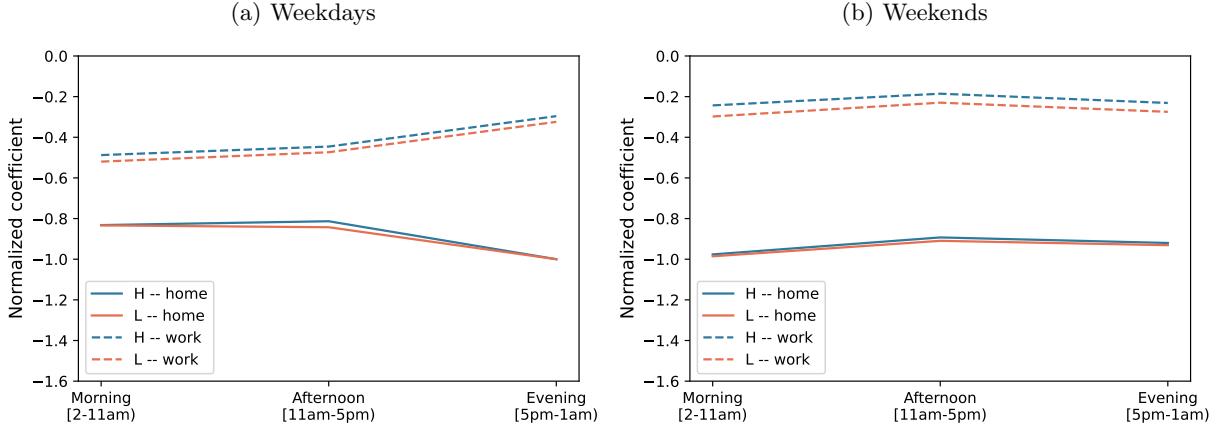
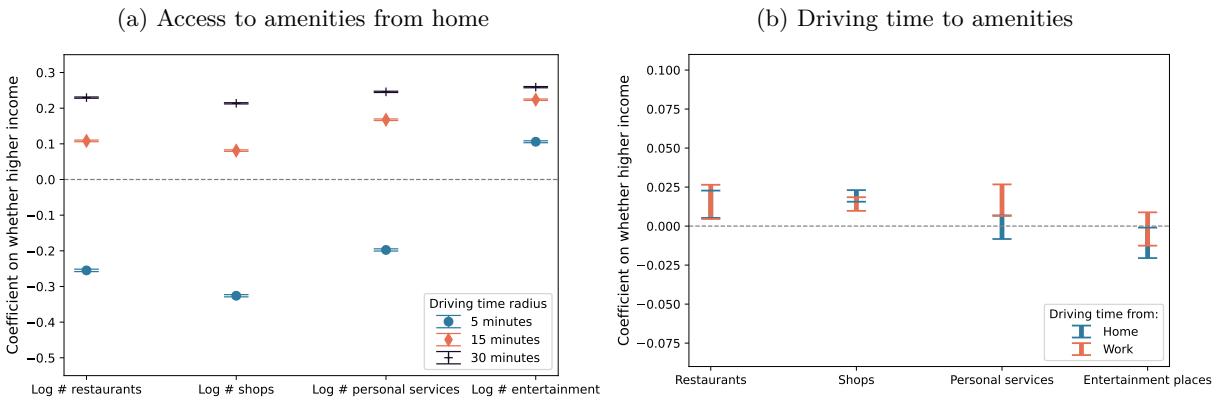
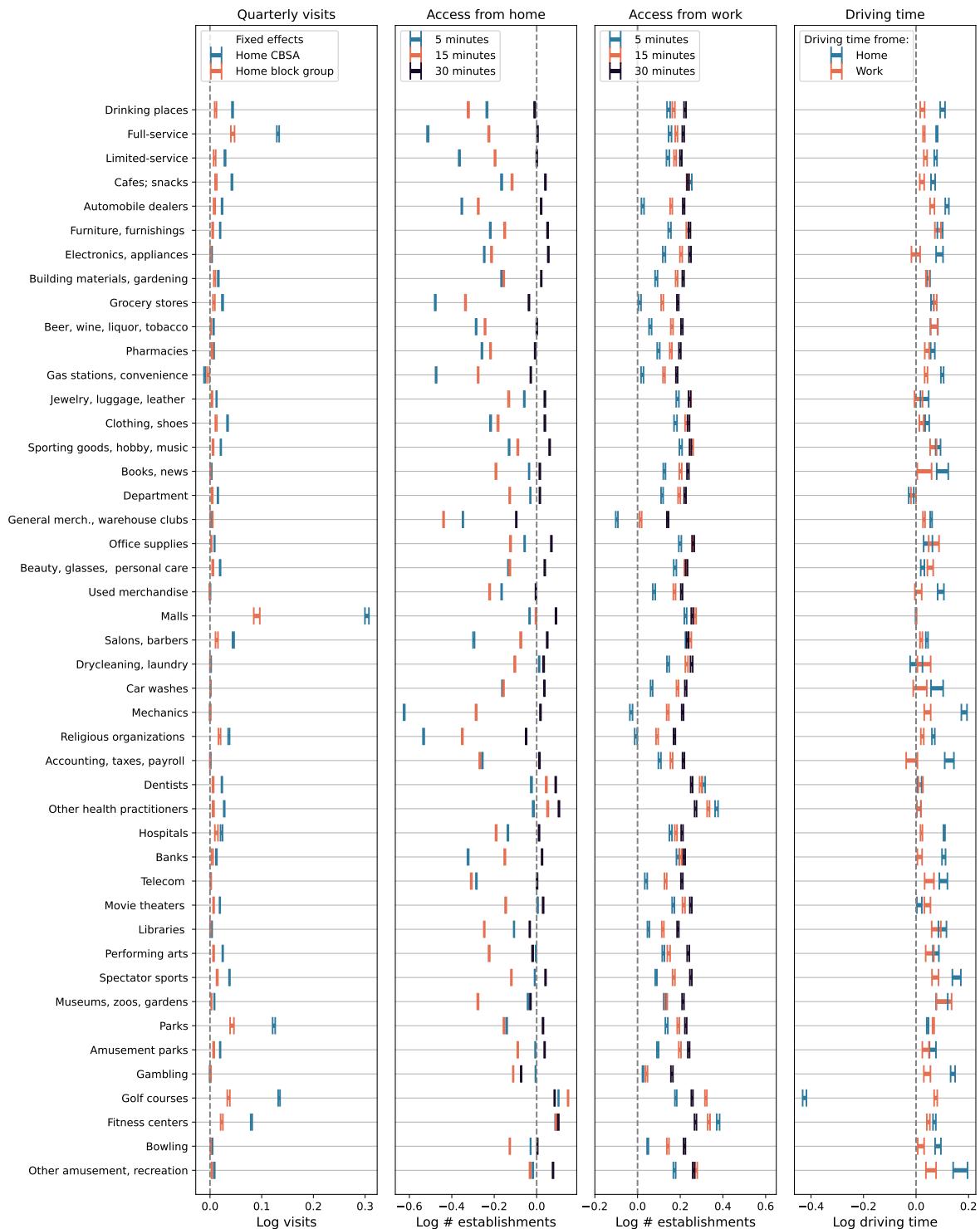


Figure C.2: Access and consumption of local amenities



*Note:* This figure plots coefficients from a series of regression where the right hand side is an indicator for whether a device is above median income. Panel a) is at the device-quarter level and regresses the log number of establishments within different lengths of driving time, controlling for the density of a device's home block group and a device's home CBSA. Panel b) is at the visit level and regresses the log of driving times from home and work on whether the device is higher income, with controls for a device's home block group.

Figure C.3: Access and consumption of local amenities (subcategories)



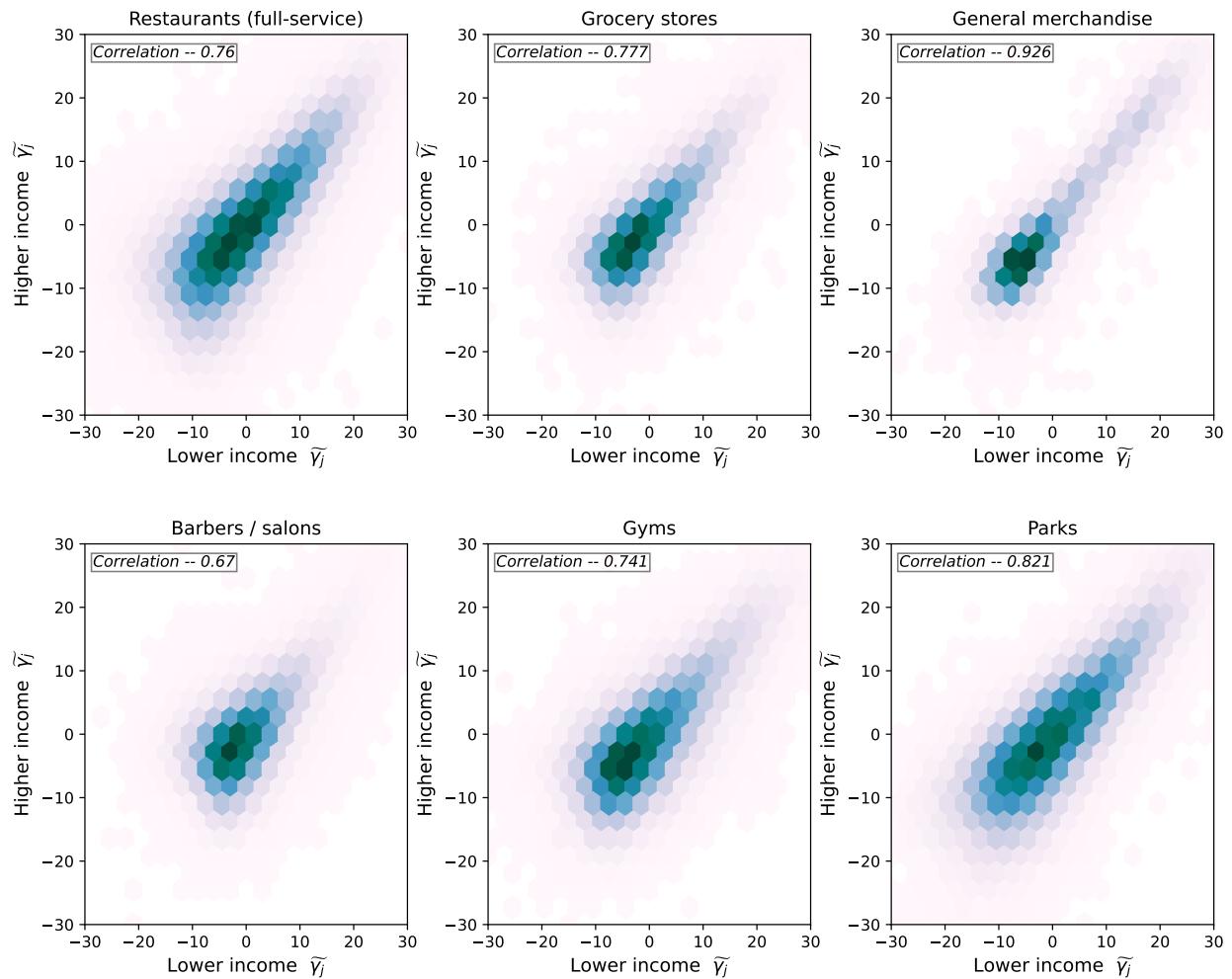
*Note:* This figure replicates Figure 1 at the subcategory level; see the note on that figure for more details.

Table C2: Top brands within each subcategory

Sub-category	Highest (H - L)	Highest (L - H)	Top H	Top L
<b>Restaurants</b>				
Limited-service	Einstein Brothers	Church's Chicken	Red Robin Gourmet Burgers	Chuck E. Cheese's
Cafes; snacks	Panera Bread	Dunkin'	Freddy's Frozen Custard	Freddy's Frozen Custard
Full-service	Zoë's Kitchen	Krystal	Texas Roadhouse	Texas Roadhouse
<b>Shops</b>				
General merch., warehouse clubs	Costco Wholesale Corp.	Family Dollar Stores	Costco Wholesale Corp.	Walmart
Grocery stores	Trader Joe's	Save-A-Lot	Meijer	Meijer
Gas stations, convenience	Costco Gasoline	Thorntons	Love's Travel Stops and Country Stores	Love's Travel Stops and Country Stores
Automobile dealers	Audi	Mitsubishi Motors	Toyota	Toyota
Furniture, furnishings	At Home	Carpet One Floor & Home	At Home	At Home
Department	Kohl's	Ross Stores	Kohl's	Kohl's
Office supplies	Hallmark Cards	Party City	Staples	Staples
Electronics, appliances	Best Buy	Cellular Sales	Best Buy	Best Buy
Building materials, gardening	Ace Hardware	Menard's	Menard's	Menard's
Clothing, shoes	DSW (Designer Shoe Warehouse)	Rainbow Shops	Marshalls	Burlington
Sporting goods, hobby, music	Dick's Sporting Goods	GameStop	Hobby Lobby Stores	Hobby Lobby Stores
<b>Personal services</b>				
Mechanics	Christian Brothers Automotive	Goodyear Auto Service Centers	Mavis Discount Tire	Safelite
Salons, barbers	Massage Envy	Regal Nails Salon & Spa	SmartStyle Family Hair Salons	SmartStyle Family Hair Salons
Banks	Charles Schwab	ACE Cash Express	First Convenience Bank (FCB)	First Convenience Bank (FCB)
Telecom	Verizon Wireless	Boost Mobile	Verizon Wireless	Verizon Wireless
<b>Entertainment</b>				
Fitness centers	Pure Barre	Planet Fitness	Lifetime Fitness	LA Fitness

*Note:* This table documents which brands have the largest gaps in the preferences of higher (H) and lower (L) income residents as well as the overall most liked brand for higher and lower income. I subset to just those categories for which there are at least 5 brands with over 100 locations and consider only those brands that have at least 100 locations.

Figure C.4: Correlation of establishment-level parameters



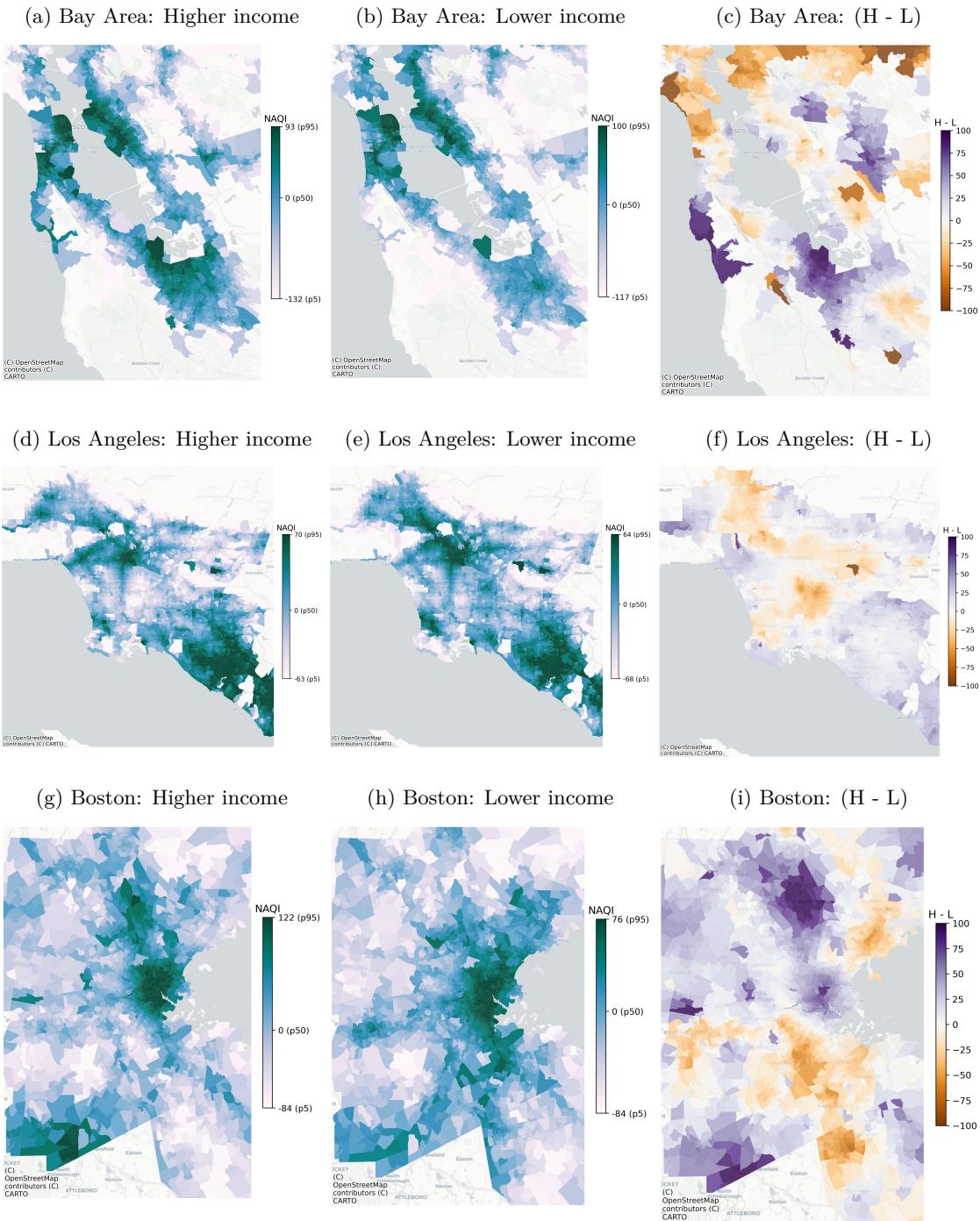
*Note:* This figure shows the joint density of higher and lower income residents' establishment-level preferences for establishments in different categories of amenities.

Table C3: Upper level estimation results

	Higher income	Lower income
<b>Arrival rates (<math>\vec{\lambda}</math>)</b>		
Weekday morning	0.0764 [0.0092]	0.072 [0.0062]
Weekday afternoon	0.121 [0.0158]	0.1048 [0.0107]
Weekday evening	0.0952 [0.0115]	0.0868 [0.0084]
Weekend morning	0.1208 [0.0128]	0.113 [0.0108]
Weekend afternoon	0.1984 [0.0219]	0.1757 [0.0167]
Weekend evening	0.1597 [0.0164]	0.1489 [0.0138]
<b>Measure of nest independence (<math>\vec{\rho}</math>)</b>		
Restaurants	0.5391 [0.2386]	0.4982 [0.2473]
Shops	0.7801 [0.1618]	0.8431 [0.153]
Personal services	0.9267 [0.0984]	0.9307 [0.1031]
Entertainment	0.9622 [0.112]	0.9526 [0.1037]

*Note:* This table the average arrival rates ( $\vec{\lambda}$ ) and upper nest measures of independence ( $\vec{\rho}$ ) for each income group. Recall when interpreting arrival rates that they hold only up to the normalization of the outside option being zero at all times of week; a higher arrival rate may imply that arrivals are truly more common or that the outside option is more valuable at that time. Each entry is an average across CBSAs, weighting by total population. Standard deviations are in brackets.

Figure C.5: Neighborhood values: Bay Area, Los Angeles, Boston



*Note:* This figure illustrates the estimated block group level neighborhood values in Equation 4.8 for a subset of CBSAs, often zoomed in on the urban cores.