

Using Machine Learning to Identify and Predict Gentrification in Nashville, Tennessee.

By

David Knorr

Thesis

Submitted to the Faculty of the Graduate
School of Vanderbilt University in partial
fulfillment of the requirements for the
degree of

MASTER OF SCIENCE

in

Earth and Environmental Science

August 9, 2019

Nashville, Tennessee

Approved:

Jonathan Gilligan, Ph.D

Janey Camp, Ph.D

Contents

	Page
List of Figures	iv
List of Tables	v
Introduction	1
Historical Trends in US Housing	1
Defining Gentrification	3
Gentrification Research	7
Theory	8
Costs and Benefits	11
Quantifying Gentrification	17
Nashville	25
Costs of Growth	26
Data Sources	27
Scope	29
Identifying Gentrification Through Machine Learning	30
K-Means Clustering Neighborhood Change	31
Nashville Neighborhood Change	34
Sensitivity	40
K-Color Join Count Statistic	41
Supervised Machine Learning	43
Random Forest	44

Prediction Performance	45
Mapped Predictions	48
Teardowns.....	51
 Discussion	 54
 Conclusion	 59
 References	 61
 Appendix.....	 68
Preprocessing.....	68
Logistic Regression	70
Logistic Regression Model.....	71

List of Figures

Figure	Page
Figure 1. Historical and projected population growth in urban areas (World Health Organization, 2018).....	2
Figure 2. Decision tree schematic of gentrification versus revitalization.....	5
Figure 3. Gentrification peer-reviewed publications (1979 – present).	7
Figure 4. Rent gap theory (Clark, 1995).	9
Figure 5. General prescriptive framework used to quantitatively identify gentrification.....	18
Figure 6. Variability in areas identified as gentrified within Davidson County, Tennessee (2000 - 2010).....	21
Figure 7. Variability between Freeman's (2005) 50 percent and 40 percent eligibility criteria for gentrification.....	22
Figure 8. Davidson County short term rental permits (2015 - 2018) (Metro Codes Department, 2019).	27
Figure 9. Davidson County home sales as a function of distance to downtown (Davidson County Tax Assessor, 2018)....	30
Figure 10. K-means clustering K-selection metrics (silhouette width and total within sum of squares).	33
Figure 11. Cluster geometries based on cluster size specifications. Plotted against the first two principle components (73.9% of total dataset variance).....	33
Figure 12. Box and whisker plots of neighborhood change typologies.	36
Figure 13. Davidson County normalized neighborhood change typologies (K=4).....	37
Figure 14. Davidson County neighborhood change cluster map (K=4)	38
Figure 15. Zoomed locations and neighborhoods of gentrification typology.	40
Figure 16. Davidson County K-means clustering (K) sensitivity (K2:K5).	41
Figure 17. Comparison of first and second-order queen's adjacency matrix and sample census tract.	42
Figure 18. Random Forest trained model Receiver Operating Characteristic (ROC) curve.....	45
Figure 19. Random forest model prediction variability. Median prediction values (x-axis) plotted against 200 model run values.	48
Figure 20. Random forest model future predictions using 2016 data.	49
Figure 21. Zoomed gentrification predictions.	50
Figure 22. South Street in the Nations neighborhood (Google Street View).	52
Figure 23. Residential teardown in the East Nashville neighborhood (Google Street View).	52
Figure 24. Gentrification predictions (2016 data) plotted against residential teardown locations and new construction values (2016 - 2019) (Metro Codes Department, 2019).....	53
Figure 25. Principle component analysis variance scree plot.....	69
Figure 26. PCA biplot.	70
Figure 27. PCA-logistic regression receiver operator characteristic (ROC) curve.	73

List of Tables

Table	Page
Table 1. Eligibility and gentrification criteria as applied by Freeman (2005), Ellen and O'Regan (2010), and McKinnish et al. (2010).	20
Table 2. Six socio-economic proxy variables used in the k-means clustering procedure (calculated as percent change 2000 - 2016).	32
Table 3. Cluster Summary Statistics (K=4).	35
Table 4. Neighborhood change typologies and associated join-count statistic p-values.....	43
Table 5. Random Forest Predictive Variables (2000).....	44
Table 6. Random forest model confusion matrix.	46
Table 7. Random forest model performance statistics.	46
Table 8. Random forest variable importance.	47
Table 9. PCA eigenvalues.	72
Table 10. PCA-logistic classification model parameters.	72
Table 11. PCA-logistic regression confusion matrix.....	73

Introduction

This research contributes to the body of academic research on gentrification by providing a data-centric methodology to identify typologies of neighborhood change. We employ an unsupervised clustering technique to group the dominant trajectories of neighborhood change in Nashville, Tennessee, measured by six important census change variables. We identify one emergent typology indicative of gentrification between 2000 and 2016 in 13% of sampled census tracts. This method bypasses the potential bias that traditional methods of identifying gentrification carry. Additionally, we argue that this method provides a more holistic and city-specific reporting of gentrification relative to other types of neighborhood change. The misrepresentation of gentrification may have important consequences for downstream analyses that often inform policy and support programs. As such, the development of gentrification as a dependent variable merits a critical examination and transparent reporting of sensitivity.

The second component of this research is to develop a predictive model to identify at-risk neighborhoods for future gentrification. This predictive component of this research is predicated on the correct identification of gentrification outlined above. We use demographic, housing, occupational, transportation, and amenity data from 2000 to train a model on baseline characteristics that may help distinguish gentrification apart from other types of neighborhood change. We report model performance for both a principle component logistic classification (PCLC)¹ and random forest (RF) models. The RF model is projected onto 2016 data to identify potentially vulnerable areas of future gentrification based on their starting characteristics.

Historical Trends in US Housing

Gentrification is a legacy of historical patterns of residential restructuring that trace back to the start of the 20th century. The end of World War II marks a mass exodus of affluent, white households from urban areas to the suburbs. The “white flight” movement was reflective of the preferences, policies, and market conditions throughout the fifties and sixties. Specific factors including systemic racism, redlining, municipal disinvestment, and the construction of the federal highway system helped to systematize the suburban movement (Zuk et al., 2015). Concurrent with this process was increased disinvestment, crime, and urban decay across many urban areas. By 1967, riots and civil unrest led to the Lyndon B. Johnson administration commissioning a report to undertake the cause of widespread racial disorder. The report concluded “Our nation is moving toward two societies, one black, one white-separate and unequal” (Kerner J.R., 1968).

¹ PCLC results are detailed in the appendix.

Economic transformations and the passing of the Fair Housing Act of the 1968 symbolize a shift from an era of institutionally racialized suburbanization towards a period of diversification and increasingly complex metropolitan structures. After 1970, the suburbanization movement was further bolstered by anti-discriminatory job and housing policies while urban cores were largely neglected. Between 1970 and 1990, two thirds of central city census tracts experienced income losses relative to their larger metropolitan areas. This trend reversed course over the 1990's with more than 40 percent of inner-city tracts experiencing increases in relative income (Ellen and O'Regan, 2010). The causes of this reversal are explored in subsequent sections, but point towards shifts in consumer preference, occupations, and the potential for profit (Ley, 1980; Smith, 1979).

During this time, gentrification attracted a significant amount of scholarly attention. The process represented a theoretical battleground between proponents of supply-side versus demand-side explanations. Although symbolically important, early accounts were largely dismissive of the idea that the process could bear the importance that it holds today. Throughout the 1980's gentrification was referred to as a *"localized small scale process...purely temporary and of little long-term significance"* (Smith, 1982). Berry (1982) similarly compares gentrification to *"islands of renewal in seas of decay"*. Others backed the idea that the process was not only part of larger, albeit infant *"back-to-the-city"* movement, but also had the potential to *"reverse the historic decline of the central and inner city, and should be actively supported by federal urban policies* (Laska and Spain, 1980).

As we have seen over the previous fifty years, gentrification is neither temporary nor is it exclusive to the United States. The phenomena has manifested in cities around the world including London, Sydney, Berlin, Hong Kong, Spain and Singapore (Atkinson, 2000; Vicario and Martinez Monje, 2003; Ye et al., 2015). The re-urbanization of post-industrial cities across the United States is a microcosm of a global trend; 68 percent of the world's population is expected to live in urban areas by 2050 (United Nations Department of Economic and Social Affairs, 2018). This has resulted in the persistent concern over gentrification on a scale today that is much larger than many of the early theorists may have anticipated.

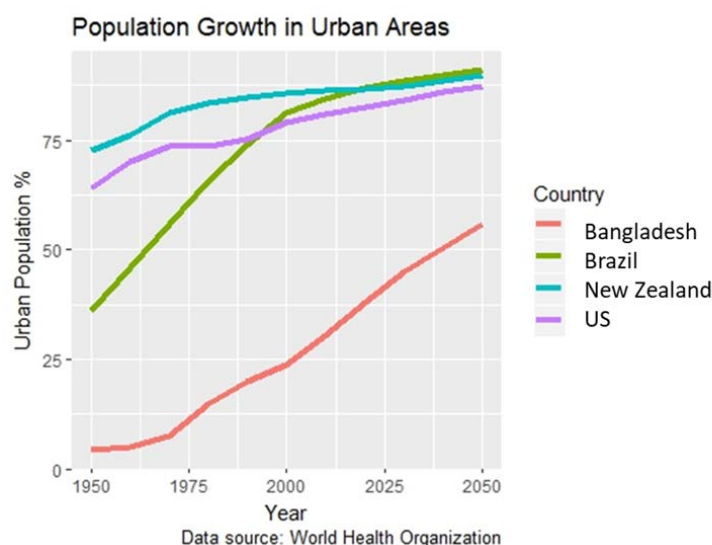


Figure 1. Historical and projected population growth in urban areas (World Health Organization, 2018)

Defining Gentrification

The conceptual definition of gentrification has been contested ever since Ruth Glass first coined the term in the 1960's. An agreed upon definition does not exist today. In this section, we emphasize the importance of defining the process in question. Gentrification has a broad, inconsistent, and colloquial public understanding. The inconsistencies in the public discourse over gentrification are largely shaped by subjective experience to the perceived threats or benefits of the process. Gentrification is a glaringly subjective process with disproportionate effects among residents. This reality is a considerable driver of not only the inconsistent public perceptions surrounding gentrification, but has also permeated and been advanced by academic research as well. The loose conceptual definitions are anything but trivial as they frame researchers' operational, or quantitative, identification of the process. These upstream aspects are central to any gentrification analysis. As such they require a critical evaluation of differences and a formal refinement for clarity.

The Department of Housing and Urban Development (HUD) defines gentrification as "a form of neighborhood change that occurs when high-income groups move to low-income areas, potentially altering the cultural and financial landscape of the original neighborhood" (United States Housing and Urban Development, 2018). The HUD definition focuses on the economic profile of gentrifiers and the incumbent populations while also broadly acknowledging the social and economic impacts of gentrification.

Alternatively, the Center for Disease Control (CDC) defines the process simply as "the transformation of neighborhoods from low value to high value" (United States Center for Disease Control, 2017). The CDC-based definition offers a one-dimensional, neighborhood-centric version focusing on the economic outcomes of the neighborhood. The definition implicitly focuses on a neighborhood's value level rather than the rate of change.

Oxford Dictionary (2019) defines gentrification as "the process of renovating and improving a house or district so that it conforms to middle-class taste." This definition solely focuses on the improvements to the built environment and overlooks any possible changes in the demographic changes that are associated with increased neighborhood investment.

Others feel the term gentrification has strayed so far from its original origin that it has been rendered unproductive (White, 2015). These dialogues have only recently been taking place, as researchers attempt to reorient the term back towards a pragmatic use. Over the course of several decades, gentrification's meaning and conceptualization gradually expanded to encompass all forms of inner-city upgrading. Adding to this de-contextualization, has been its emergence on the global stage, which has brought unlimited diversity in the processes that shape neighborhoods and their associated characters, responsible actors, and spatial limits. Maloutas (2012a) argues that the definition has been stretched unrecognizably beyond its original intent, resulting in "a regression in conceptual clarity and hence in theoretical rigor." Halle and Tiso (2014) contend that gentrification is used "“very loosely,

conflating several issues that should be considered separately.” Similarly, Stern and Seifert’s (2007) rationalize that “..if we see neighborhood revitalization as desirable, we cannot afford to label all population change as gentrification.” Atkinson (2004) states “.. the term gentrification is predicated on displacement and community conflict.”

The examples given above provide a glimpse at how the vocabulary surrounding gentrification varies considerably. The disagreements between definitions go beyond semantics. They frame how we perceive and discuss the benefits, costs, causes, and ultimately the responses to gentrification. Varying conceptual definitions provide an unstable foundation for how we measure and identify gentrification, in practice. The lack of consistency within the language of gentrification has unsurprisingly driven inconsistent methods used to measure the process. The inconsistent methods spur conflicting results which further obscure our understanding of gentrification in a non-constructive feedback loop. We propose a re-examination into the methods that researchers use to measure and identify gentrification. Both accurate conceptual and operational definitions of gentrification are foundational to our efforts towards creating more socially just cities.

Any improvement on our conceptual definition of gentrification must acknowledge the context and specific effects of the problem. It is therefore beneficial to frame gentrification as one specific type of neighborhood change. We consider neighborhood change as an ongoing spatial and temporal process associated with the movement of people and capital. Gentrification is distinguished by a sustained period of disinvestment, followed by an influx of investment and wealthier residents that results in the displacement of existing residents. This definition provides a much clearer pathway towards operationalization by suggesting distinguishable and measureable signals rather than more subjective ideals like “middle-class taste.” Furthermore, this definition includes the displacement of incumbent residents as a prerequisite feature of gentrification. The causal link between gentrification and displacement has been the target of an entire subfield of literature, but can be bypassed by refining the term to require displacement as one of its most critical and foundational attributes. Displacement is the critical negative externality that differentiates the process apart from other types of neighborhood change such as revitalization and/or incumbent upgrading. Efforts to separate gentrification from displacement may be viewed as potential scrubbing of the most important defining consequences of the process. Thus, we will refer to gentrification without displacement as revitalization (Clay, 1979; Zuk et al., 2015).

Revitalization, also called incumbent upgrading, is considered a positive outcome in which public and private investments work collaboratively with community partnerships and local advocacy groups to build opportunity and better living conditions for those residents living in previously disinvested communities (Tatian et al., 2016). Revitalization fundamentally retains the demographic, socioeconomic, and cultural composition of a neighborhood (Helms, 2003). Revitalization strategies work to combat displacement by ensuring existing resident’s capacity to stay in their homes and realize the benefits of neighborhood improvements. Specific strategies attempt to preserve affordable housing, encourage affordable housing developments, and engage existing community residents (United States Housing and Urban Development, 2018).

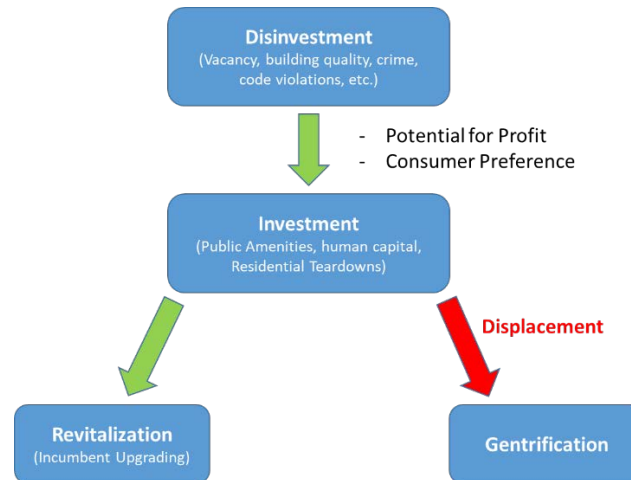


Figure 2. Decision tree schematic of gentrification versus revitalization.

Our conceptual framework for gentrification builds off of Lees (2007) comment: “We need a more fine grained approach in gentrification research, one that is both more specific and more general”. The working definition presented here is both more specific by requiring the displacement of existing residents, and more general, in that it attributes the social and cultural changes to a broad class of measureable investments. Displacement may be quantitatively captured through population flows—specifically race, but also through class indications such as income levels and educational attainment. Likewise, investments and disinvestments are reflected in home/rent values, human capital (education), public amenities, businesses, and various other physical or social signals. Our effort to quantitatively identify gentrification builds off of these multi-dimensional and often overlapping proxies to examine neighborhood change more holistically.

Displacement

We argue that displacement is a fundamental aspect of gentrification and one that merits further refinement. Previous researchers have discussed the causal connection between gentrification and displacement, but fall short of requiring displacement as a constitutive element and defining aspect of the gentrification process. In this section we define and explore the causal evidence for gentrification and displacement that have served more to obscure than to clarify. These inconsistent effects are a consistent theme of gentrification research that have their root in quantitative issues discussed in subsequent sections.

George and Eunice Grier (1978) present an early definition of displacement: “*Displacement occurs when any household is forced to move from its residence by conditions that affect the dwelling or its immediate surroundings, and that: 1) are beyond the household’s reasonable ability to control or prevent; 2) occur despite the household’s having met all previously imposed conditions of occupancy; and 3) make continued occupancy by that household impossible, hazardous, or unaffordable.*” Marcuse

(1986) identified several overlapping modes of displacement that are manifested through physical, economic, chained, or exclusionary mechanisms. Marcuse argues that efforts to understand the link between gentrification and displacement may be under representative because of significant time lapses between abandonment of properties and gentrification. As a result, “chains” of displacement as well as exclusionary displacement must be taken into account, but they are often invisible to aggregate data sources and go uncounted.

An entire subfield of literature has focused on the link between gentrification and displacement. While there are conceptual and methodological differences in identifying gentrification alone, a consensus for measuring displacement may be even more varied. Displacement presents analytical challenges that aggregated data like the US Census are ill-suited for considering the longitudinal data needed to track displacees as well as their reasons for moving. Atkinson (2000) likens measuring displacement to measuring the invisible. Distinguishing gentrification-induced displacement apart from voluntary migration and incumbent upgrading make measures of displacement increasingly more difficult (Atkinson, 2000). As a result, a majority of studies attempting to quantify the magnitude of displacement employ a narrow definition covering only evictions or unaffordable price increases (Zuk et al., 2015). The inability to consistently and effectively measure displacement has stalled local and federal government efforts to move from gentrification to revitalization.

A study by the Philadelphia Federal Reserve Bank concluded that low-income residents were no more likely to move from their homes in a gentrifying neighborhood than a non-gentrifying one (Ding et al., 2016). Lance Freeman’s 2005 study concluded the potential for household to be displaced in a gentrifying neighborhood to be only 1.3% on a national scale (Freeman, 2005).² In another national study, Vigdor (2002) similarly found no evidence that gentrification increased the probability that low-income families exit their households. Atkinson (2000) found above average losses for socio-economically vulnerable groups during the 1980’s in London suggestive of working-class displacement. These studies grabbed considerable media attention with Time magazine’s headline “Gentrification: Not Ousting the Poor” (Kiviat, 2008; Wyly et al., 2010).

Conversely, Newman and Wyly (2006) re-examined Freeman’s displacement data for New York City data and found that 6-10% of all moves between 1989 and 2002 were attributed to some form displacement. New York City has been a test bed for displacement research largely because of the New York City Housing and Vacancy Study (NYCHVS), a longitudinal survey that captures renter moves and more importantly, their reasons for moving (Wyly et al., 2010). 17 percent of African American poor renter moves between 2005 and 2008 were attributed to forced displacement. Additionally, the median age for renter moves due to displacement was 46 compared to 35 for all other move reasons. The researchers concluded an estimate of approximately 10,000 displacement-induced moves per year in New York City alone (Wyly et al., 2010).

² The results of this study are national, likely masking significant heterogeneity between cities.

Gentrification Research

British sociologist, Ruth Glass, first coined the term gentrification in the late 1960's, describing the process she observed in London, by which poor, inner-city ghettos were displaced by increasingly affluent neighborhoods (Glass, 1964). Her influential work, *London: Aspects of Change* (1964), drew upon academic fields of sociology, geography, history, medicine, and urban planning. The phenomena she observed in London would later be dissected by American scholars in the decades following in an effort to advance the theory, causes, history, and responses to gentrification. At the time of Glass' publication, the nascent views of housing economics rested on a filtering and succession model. This theory suggested that newer, more spacious suburban housing is filled by well off families looking to upgrade while their previous smaller home is filled in by poorer occupants (McKenzie et al., 1925). Gentrification existed within this model only as an exception to the rule and necessitated a more complex explanation than the filtering model could afford. In the following section, we outline the distinct trajectory that gentrification research has evolved. Theoretical focus of the causes of the process were advanced through two decades of disagreement between proponents of the supply-side versus demand-side explanations for gentrification. The empirical research on gentrification that emerged since the 1990's has largely focused on characterizing the consequences of gentrification.

The problem of gentrification has not been hindered by a lack of scholarship interest. Publication trends of gentrification-related studies in the peer-reviewed literature (Figure 3) show a rapid and pronounced rise over the last two decades, which cuts across academic disciplines such as urban studies, geography, environmental studies, sociology, economics, history, and political science (Clarivate Analytics, 2019).

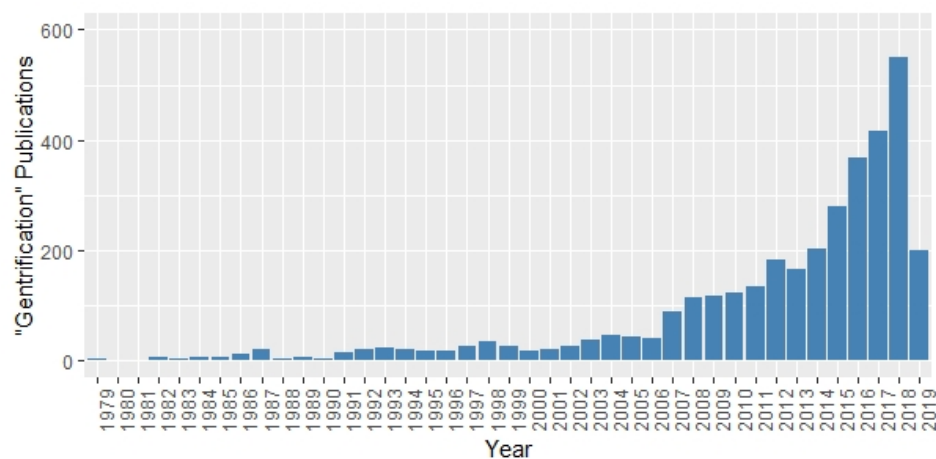


Figure 3. Gentrification peer-reviewed publications (1979 – present).

Theory

Early studies of gentrification were centrally concerned with appointing the causes underlying the process that had, at the time, been limited to larger cities such as London, New York, and San Francisco. Two competing theoretical movements waged a high stakes battle on the topic throughout the 1980's. David Hamnet (1991) offers an explanation for the theoretical discord writing *"...Gentrification is a frontier not just physically, economically, socially and culturally, but also theoretically, ideologically, and politically. It comprises a contested boundary zone between radically different theories and explanations. And it is arguably this aspect of gentrification, above all others, which has kept the gentrification debate at the forefront of urban geographical literature for over a decade."* The boundary zone he described is differentiated by *"proponents of the imperatives of capital and profitability"* on one side and *"the proponents of culture, preference, and human agency"* on the alternate side. This theoretical divide is alternatively expressed as structural Marxism versus liberal humanism, supply-side versus consumption-side economics, or by the flow of capital versus the flow of people.

Supply Side

The first movement advancing its theory of gentrification contended that the process is born out of the macro-economic patterns of investment and disinvestment in the built environment. This theory was put forth by Neil Smith as early as 1979. It describes the inertia of the macro-economy that precipitates the conditions necessary for gentrification to occur. Smith established the rent-gap theory to explain investment decisions triggered by a disparity between the current rental value of a property and its potential rental value (Smith, 1979). Figure 4 shows building values declining over time due to physical deterioration, under maintenance, and style obsolescence. Capitalized land rent reflects the value of land under its current designated use, and value obtained from the current structure, services, and/or improvements built on it (Smith, 1987). Sales prices are a function of both building value and capitalized land rent. The final rent gap factor is potential land rent, or the unrealized land rent that could be achieved under a more optimal (redeveloped) land use. Prolonged disinvestment drive sales prices down until there is a sufficient rent gap to trigger reinvestment. Smith explains that as a neighborhood declines, the rent gap widens until redevelopment becomes a rational market response. It is only at this point that redevelopment makes economic sense as there would be unjustifiably little profit from redevelopment in areas that already maximized their current land rent (Smith, 1979).

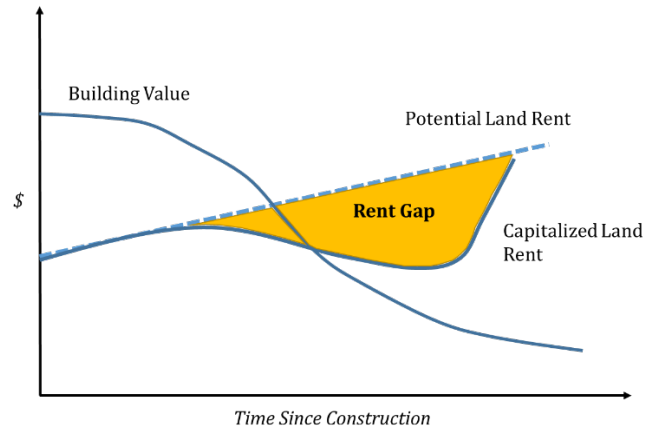


Figure 4. Rent gap theory (Clark, 1995).

Smith's theory also offered insight into the spatial distributions of investment and disinvestment. He emphasized that uneven development and investment patterns are inherent to capitalist economies as rational actors search for specific areas that will maximize profits. Investment in one area creates barriers to further development in that same area. Likewise, focused investment in that area leads to underdevelopment of other areas due to the zero-sum nature of capital investment. This theory helps make sense of the widespread suburbanization movement following World War II that was accompanied by disinvested urban cores. Additionally, this offers an explanation for the resurgence of inner-city investments that were primed by decades of neglect and represented relatively low barriers to entry. Gentrification could be seen as part of a larger economic process in which suburban housing markets saturate, profit margins in that locale shrink and underdeveloped areas in inner-cities become more attractive investment opportunities.

The rent gap theory has held up considerably to empirical research focused on explaining investment decisions in urban areas. Rosenthal and Hensley (1994) found evidence in Vancouver that housing redevelopment of vacant land is triggered when its value exceeded that of its current land use. Munneke (1996) similarly found a higher likelihood of redevelopment for industrial and commercial land in Chicago when the redeveloped value exceeds that of its current use. These intuitive findings lend support to Smith's theory that redevelopment is a rational and potentially predictable market response to changes in land use and the potential for profit.

Smith prescribed consumer preferences as an important, but secondary cause of gentrification, declaring *"The so-called urban renaissance has been stimulated more by economic than cultural forces. In the decision to rehabilitate inner city structure, one consumer preference tends to stand out above the others - the preference for profit, or, more accurately a sound financial investment."* (Smith, 1979). Later academics would dispute Smith's weighting of the forces behind gentrification, arguing for equal or higher standing of individual actors, policies, and preferences.

Consumption Side

The alternate theory emphasized consumer preferences, human agency, and consumption side demand as the necessary forces for gentrification. David Ley (1980) argued that the flow of people, as opposed to the flow of capital, is the largest force behind of gentrification. In his paper, *“Liberal Ideology and the Post-Industrial City”*, Ley examined the critical elements that characterize a shift from industrial cities to post-industrial cities. Among these factors are the labor force transition from blue collar jobs to white collar jobs, manufacturing to service-based urban economies, and rapid technological advancements. Ley’s post-industrial thesis highlights the role of individual actors’ lifestyle preferences, values, and aesthetics. In this model, it is the combination of labor market transitions and changing values that produces a sufficient pool of potential gentrifiers. Consumption-side efforts to explain gentrification therefore begin with the individual actor’s demand for inner-city amenities. Ley and others prioritize this population’s ability to “restructure the built environment and accelerate gentrification” (Ley, 1986).

Richard Florida (2003) distinguished the important factors that have transformed the economic and cultural landscape of today’s city. Florida’s thesis focuses on the consequences of the shift away from a material-based economy towards a creative, skill-driven, and human-capital-powered economy. The “creative class”, he argues, is “the key force that is reshaping our geography, spearheading the movement back from outlying areas to urban centers” (Florida 2012). The creative class comprises nearly 40 million American workers, all fundamentally working to “create meaningful new forms” (Florida, 2003). The creative class is distinguished by occupations ranging from artists and professors to software engineers, and finance professionals. This class is associated with shifting values that emphasize “belonging, self-expression, opportunity, environmental quality, diversity, community, and quality of life” (Florida, 2014). Florida argues that these group’s values in combination with economic shifts results in the renewed demand to live in cities as creative centers.

Contemporary Theory

Over the course of the past 50 years, our perception of gentrification has evolved from a sporadic, localized process to a pervasive threat to disadvantaged communities across the globe. The theoretical discord throughout the 1970’s and 1980’s proved beneficial to refine the governing principles of the process that are deeply rooted in profitability, post-industrial economic transformations, and human agency. This research era was fundamentally concerned with the “why” of gentrification. The process conflicted with the traditional filtering models that assumed continual urban decline and the logical response towards suburbanization. Gentrification necessitated and received a nuanced explanation of urban upgrading that highlighted social and economic factors. This early research was largely beneficial towards providing a clear avenue to ask more specifically “how” and “where” gentrification may manifest in certain cities. Despite this research narrative, we will see that gentrification literature largely jumped towards studying the causes and effects before an agreed upon conceptual and operational definition had been agreed upon.

Both consumption and supply-side explanations offer much to a comprehensive understanding of the causes of gentrification. To summarize, supply-side proponents attribute gentrification to

disparities between housing value and the untapped profit that could be achieved under an alternate land use. Consumption-side explanations suggest economic shift to a skills-based economy have produced young, middle-class professionals who value the diversity and amenities that cities offer.

Critical to each theory is the starting point to explain conditions necessary for gentrification. Each theory offers a necessary, but incomplete explanation of gentrification. Smith's rent-gap explanation is most useful when framed not as a deterministic trigger for gentrification, but rather as a catalyst that can provide valuable insight into which areas may be primed for potential reinvestment. Of primary importance is Smith's prerequisite for a sustained period of disinvestment prior to gentrification. Disinvestment is a measureable process through various proxies that include vacancy rates, housing construction age, as well as home and rent values. The rationality behind actor's investment decisions lends confidence to the primary assumption of this research- that gentrification is structured and to some degree, a predictable process.

However, without the existence of potential gentrifiers, acting on their preference for inner-city living, gentrification does not exist, irrespective of the rent gap. The dueling explanations for gentrification pioneered by Smith and Ley have been refined by other researchers towards a more complimentary and comprehensive explanation for the causes of gentrification. David Hammet (1991) synthesized the polarizing theories of gentrification. He concludes that gentrification doesn't rely on a singular cause, but may emerge when there is a pool of gentrifiers with a cultural preference for urban living as well as a sufficient supply of inner-city housing. These principles are critical in order to refine the theory behind our efforts to identify the patterns of gentrification while also providing the basis for our gentrification predictions.

Costs and Benefits

The subsequent era of research on gentrification throughout the 1990's and 2000's largely shifted towards understanding the costs and benefits associated with the process. In a similar vein to the previous debate concerning the causes of gentrification, its consequences would again polarize researchers, policy-makers, and public opinion. Some viewed gentrification as a beneficial force that could stimulate local tax revenues, restore the built environment, and decrease crime. Supporters of this view contend that gentrification should be actively promoted by government policies. Others pointed to the process' social costs, citing unaffordable housing costs, industrial job displacement, and cultural cleansing. The logical response for this camp was government policy responses that could combat the market forces in order to preserve affordable housing. Here we discuss the key (and contradictory) literature that continues to fuel and reinforce each group's perceptions of gentrification.

Public Amenities

Both private and public actors may shape, but also respond to gentrification. Gentrification may trigger long overdue improvements to public spaces and local services like parks and public transit. These public improvements are often made possible by the increased tax revenues for municipalities (Lang, 1986). Although these public investments may be intended to serve the original residents, the new amenities are priced into the housing and rental markets; this can unintentionally price existing residents out of the local market before they can realize the benefits.

A substantial amount of recent research has looked at the impacts of public transportation investments. These public investments are designed to improve accessibility for low-income residents and their possible unintended consequences in raising housing prices and triggering displacement. This renewed public investment has the potential to attract wealthier residents, leading to potential displacement of the most transit-dependent riders away from transit infrastructure and towards the suburbs where transit options are sparse. Much of the academic research in this area has focused on the appreciable housing price impact of transit oriented developments (TODs). On the high end of the spectrum, housing near transit premiums were found as high as 45 percent in Santa Clara County, California (Cervero and Duncan, 2004), 24 percent in a national study spanning 12 metros (Pollack et al., 2010) and 17 percent in Chicago (McDonald and Osuji, 1995). Several more studies found appreciation rates between 5 and 15 percent including Boston (6.7 percent) (Armstrong Jr, 1994), San Francisco (10-15 percent) (Cervero and Landis, 1997), San Diego (17 percent for condominiums, 6 percent for single family homes) (Duncan, 2008), London (9.3 percent) (Gibbons and Machin, 2005), and Portland (10 percent) (Knaap et al., 2001). However, other studies found little effect or a negative effect between proximity to rail transit and property values in Atlanta, Buffalo, and Southern New Jersey (Bowes and Ihlanfeldt, 2001; Chatman et al., 2012; Hess and Almeida, 2007).

The disagreement between studies could either be reflective of methodological differences between studies or larger city-specific differences in housing tenure, extent and quality of transit systems, housing market conditions, and surrounding developments (Wardrip, 2011). Furthermore, these studies account only for housing values and do not detail the link between increased home values and displacement.

Employment Opportunities

Meltzer and Ghorbani (2017) assessed the link between gentrifying low-income neighborhoods and employment opportunities- an often cited benefit of gentrification (Freeman, 2005; Vigdor et al., 2002). Their results found that on average, gentrifying neighborhoods do not see a meaningful increase in employment compared to non-gentrifying low-income census tracts. The question of employment opportunities in gentrifying areas is confounded by observed growth in retail opportunities at the expense of industrial jobs (Curran, 2004; Monroe Sullivan and Shaw, 2011).

A 20-city study by Lester and Hartley (Lester and Hartley, 2013) examined the effects of gentrification on job sector growth between a contemporaneous period (1990-2000) and a long-run period (1990-2008). They found an average loss of roughly three jobs per gentrifying census tract during the contemporaneous period. The bulk of jobs lost were from manufacturing industries (-28.6 average jobs), however this was countered by a significant increase (+10.1) in the restaurant industry as well as

large, but statistically insignificant gains in the service sector (+19.5). Interestingly, the researchers found a statistically significant decline in retail jobs within gentrifying tracts compared to all other tracts. Relative to only low-income tracts, however shows a net job growth of approximately 11 jobs. Additionally the researchers found much larger employment benefits using their longer term period, with 60 additional jobs relative to non-gentrifying areas and 100 jobs relative to other low-income, non-gentrifying tracts.³

Education

Increased neighborhood investment certainly has the potential to increase educational opportunities for children through sustaining high-quality public schools that are appealing to parents across different income levels (Formoso et al., 2010). Mixed income schools have been shown to benefit children from low-income families (Black 1996), however the arrival of more affluent residents may inadvertently reduce access to institutional resources through increased competition (Formoso et al., 2010). Many of the affluent in-movers to gentrifying areas are comprised of younger residents without children⁴ and which can result in the decreased enrollment in public schools. Additionally, more affluent residents may choose private schools further diminishing the enrollment rate at public schools. Researchers found an 18% decrease in public elementary school's enrollment rate in gentrifying areas of Chicago compared to a city-wide increase of 13% (Catalyst 2003).

Displacement is a destabilizing force that has significant impact in declining school performance (Phillips et al., 2015). Children in poor families are significantly more likely to move schools. Frequent family relocations are associated with a 35% increased likelihood to repeat a grade and 77% higher chance of reported behavioral problems (Wood et al., 1993).

Public Health

The CDC refers to gentrification as a public health issue, citing the lack of affordable healthy housing, healthy food options, transportation, health services, and social networks as detrimental to resident's health (United States Center for Disease Control, 2017). Gentrification has the potential to bring access to new resources for a community that may improve health outcomes. On the other hand, gentrification has the potential to displace existing residents, before the institutional benefits are realized.

For those that stay, rising rents can place a significant financial burden on existing residents, leading them to sacrifice basic needs such as health care and healthy foods (Phillips et al., 2015). Tenants may also feel disenfranchised to voice concerns over unsafe living conditions to landlords out of threat of evictions (Jonathan Gilligan, 2019). These residents are also likely to face disruptions to their

³ The researchers equate the conflicting magnitude of their results to the additional 8 years between their contemporaneous and long-run periods, but it is highly likely that broader economic patterns like a 31-year low (2000; 3.9%) and 15-year high (2008; 7.3%) could account for the inter-period differences.

⁴ The proportion of non-family households increased from 62% in 1970 to 71% in 2000 (Birch 2005).

social networks that can lead to chronic stress, increased susceptibility to long term disease and higher mortality rates (Atkinson, 2000; Mani et al., 2013; McEwen, 1998).

Hurnh and Marokor (2014) examined the correlation between gentrifying areas of New York City and pre-term birth outcomes. They found high levels of gentrification to be adversely attributed to pre-term birth for people of color, and a beneficial relationship to non-Hispanic Whites. Gibbons and Barton (2016) found a modest increase in self-reported health responses within gentrifying neighborhoods in Philadelphia. However, the benefits were distributed unequally, with nonwhite residents reporting lower self-rated health. Overall, the research on gentrification's impact on health outcomes offers a familiar mixed bag of evidence while suffering from the inability to discriminate between existing and new residents.

Hazardous waste sites represent negative community amenities and a danger to public health. To protect against these harms, the Environmental Protection Agency (EPA) places the most dangerous sites on the National Priorities List (NPL) for extensive remediation, deemed Superfund sites. Many who live near these contaminated areas lack the residential mobility to move to safer locations prior to the government sponsored cleanups. More than 15 million people, or 5 percent of the US population lives within 1 mile of an EPA-designated Superfund site. This population is comprised of 49.3 percent minorities compared to the national percentage of 38.4 percent. They are also disproportionately likely to be in poverty (16.7%), linguistically isolated (8.4%), and have less than a high school education (16.3%) (US EPA, 2017).

The location of some of the most vulnerable populations around contaminated areas stems from a larger environmental justice issue. However, their clean up and conversion into productive land may trigger gentrification. The median value of homes within one kilometer of a site appreciated by 18.5% following environmental cleanup (Gamper-Rabindran et al., 2011). Perhaps even more telling are the demographic changes associated with the government funded cleanups. Gamper-Rabindran and Timmins (2011) found EPA remediation of hazardous sites on the NPL led to a 26 percent increase in mean household income and a 31 percent increase in college graduates.

Crime

Public safety is an important neighborhood amenity potentially affected by gentrification. However, the link between gentrification and crime remains unclear; McDonald (1986) found decreasing personal crime rates within gentrifying areas. Contrarily, Taylor and Covington (1989) reported increases in personal crime rates and decreases in property crime within gentrifying areas. Kreager et al. (2011) found the inverse; short term increases in property crime and no association between violent crime. Linking specific crimes, like assault, to gentrification has resulted in similar contradictory results; Lee (2010) and Van Wilsem et al. (2006) find positive associations, while O'Sullivan (2005) reports a negative association. Barton (2014) found sub-boroughs of New York City that experienced gentrification had significant decreases in assaults, homicides, and robberies.

Renter-Homeowner Dichotomy

Homeowners and renters may have very different outcomes in the face of gentrification and displacement. Homeowners are often able to capitalize on rising property values and receive a positive return on investment. However, it is important to highlight the financial inequality between renters and homeowners when discussing the costs associated with gentrification. Firstly, the financial importance of homeownership cannot be understated; the median net worth of homeowners in the United States (\$231,400) is 46 times that of renter households (\$5,000) (Joint Center for Housing Studies, 2018). As such, the pathway to displacement for these separate groups looks different in the face of gentrification. Renters are often burdened with a more direct financial pressure when a landlord decides to raise rent, oftentimes with little notice or justification.

For homeowners, rising property taxes are often considered the primary mechanism leading to financial pressures and displacement. A Philadelphia-based study found substantial increases in property assessments and property taxes within gentrifying neighborhoods. The researchers found that gentrification led to an average property increase of over \$70,000, equating to a \$542 yearly tax increase (including tax relief programs) (Ding et al., 2016). Displacement for homeowners is a more indirect process than that of renters. Homeowner displacement requires spillover from adjacent development to home values, property assessments, and ultimately increased tax rates as the ultimate mechanism for displacement. However indirect, the increased taxes can be especially burdensome specific subpopulations, like elderly homeowners operating on fixed incomes. The median age for renter moves due to displacement in New York City was found to be 11 years older than relocation due to voluntary moves (Wyly et al., 2010). In Philadelphia, however, researchers found that elderly residents represent some of the least likely groups to move (Ding and Hwang, 2016). Elderly homeowners in Philadelphia may be able to burden the tax increases because of local tax relief programs which may not be available to elderly residents in other metros.

Martin and Beck (2018) conducted a national study focusing on the difference between renter and homeowner susceptibilities to displacement. Their findings suggest renters are nearly twice as likely (2.6%) to involuntarily move from a gentrifying neighborhood than a homeowner (1.3%). Additionally, they found that homeowners within a gentrifying neighborhood are no more likely to move than homeowners in non-gentrifying neighborhoods (Martin and Beck, 2018). This study suggests that the mechanisms leading to displacement, whether rent hikes or property tax increases work at an uneven pace for renters compared to homeowners. Lang (1982) argues that the costs and benefits associated with gentrification are highly subjective and depend largely on the perspective of the stakeholder involved. The renter-homeowner dichotomy is one of the most important stakeholder perspectives that make it exceedingly difficult to generalize the costs and benefits of gentrification.

Cost-Benefit Conclusions

The empirical research on the consequences of gentrification clearly offer a mixed bag of findings. Gentrification can either have a positive, negative, or neutral relationship on these important outcomes. Although we only cover a portion of the literature, it is evident that not only the magnitudes, but the overall directionality of these links can be contradictory. These consequential findings have provided quantitative evidence for both supporters and critics of gentrification to further entrench their views. Implicit in these studies is the question: "Is gentrification a battle worth fighting?" The

contradictory empirical evidence presented here on the costs and benefits of gentrification makes it acutely difficult to generalize the process as wholly good or bad. It also sheds light on the careful balance that public actions like expanded transit and environmental cleanups must consider in order to serve the most disadvantaged communities without displacing them.

To summarize the consequences of gentrification, we must return to our conceptual definition of gentrification, which requires the displacement of residents. Residents who are displaced lack the resources to overcome the financial burdens that are one of the few consistent outcomes of gentrification. Therefore, it is important to highlight that those displaced early in the gentrification process will not be able to realize *any* of its reputed benefits. This arrives us at the fundamental problem with gentrification- the unequal distribution of its costs and benefits. The costs are disproportionately burdened by the most vulnerable populations in need of assistance the most. These costs are placed squarely on groups that have borne centuries worth of costs, whether in the form of slavery, segregation, racism, or discrimination. In this light, gentrification is ill-suited for analyses of its net effects. The wide range of subjective experiences and individual perceptions offer a largely unaccounted explanation to the polarizing public sentiment on gentrification.

The fractures in gentrification literature may also be accounted for by important methodological differences. Studies are able to arrive at drastically different conclusions regarding the effects, spatial extent, and causes of gentrification, even after considering subgroups and their disproportionate effects. We argue that the conflicting empirical findings are not so surprising when the methodological differences are taken into account. Different time periods, study areas, data sources, and gentrification proxies all feed into the conflicting results. There is also little attention paid to the geographic scale and contextual differences across studies (Lees, 2000). National scale studies are often limited to census data and may hide important inter-city variations. On the other hand, studies that look at individual cities may be able to tap local-level data, but likely have different demographics, policies, and external factors that only permit crude generalizations. Time scales used in analyses can also bias the impacts of gentrification; limiting the temporal scope of analysis could discount early gentrified areas. As gentrification has matured and spread, many of the displacement studies of the early 2000's may also be severely outdated.

The type of study and questions asked have significant implications for the conclusions that it may draw. Qualitative studies are generally undertaken in cities of neighborhoods that are believed to have experienced or are experiencing gentrification (Brown-Saracino, 2017). Their preselection as gentrifying areas allows researchers to ask detailed questions about how gentrification manifests and what its consequences are. Because of their narrow scope, they focus on ground level interactions between individual actors, businesses, or institutions. What they gain in detail, however, is often lost in generalizability across broader patterns due to differences in local policies, demographics, and the intensity of gentrification. On the other hand, quantitative studies generally ask questions pertaining to which areas of a city are gentrifying, at what intensity, and with what effects. Quantitative studies have a wider scope, typically using aggregate data for the city or metropolitan statistical area. While this larger scope allows researchers to evaluate the extent of gentrification relative to non-gentrifying areas, it can overlook important distinctions such as pockets of gentrifying areas within a census tract that as a whole does not appear to gentrify. Their relative view of gentrification from a bird's eye view may draw very

different conclusions about the same gentrification process within a city than their qualitative counterparts. Quantitative studies which are more removed from the process generally paint gentrification in less extreme terms, while qualitative studies highlight the emotional consequences of the process (Brown-Saracino, 2017)

Despite these glaring inconsistencies, what all gentrification studies do have in common is a method to distinguish gentrifying areas apart from non-gentrifying areas. A critical examination of these methods uncovers more inconsistency, confirmation bias, and arbitrary sensitivity. The majority of studies that measure the effects and extent of gentrification employ inconsistent, prescriptive operational definitions to identify gentrification. The poor conceptual foundation for gentrification, discussed earlier, has permitted a wide lens of operational definitions. As a result, gentrification could be measured in many different ways. Within the same study area and time period, the extent of gentrification can vary substantially, depending on the method used to measure it (Enterprise Community Partners, 2019). In the following section we explore the operational definitions of gentrification and how they may explain a considerable amount of variance found within the gentrification literature.

Quantifying Gentrification

Much of the disagreement across academic literature may be attributed to the methodological differences used to measure gentrification. Researchers have largely charged past certain red flags when translating the conceptual understanding to a measurable process. An early warning stated: "If one wants to better understand, predict and even alter changes in urban neighborhoods, one thus must be exceedingly careful in operationally specifying the exact dynamic in question, and must recognize that such a specification may, in itself, influence the outcome of the analysis" (Galster and Peacock, 1986). The warning fell on deaf ears, however, as researchers continue to vary the operational definitions of gentrification. Until recently, this discrepancy received little attention. Maloutas (2012) notes that researchers "would have plenty to gain from an increased awareness of the contextual limits of their own tools." Similarly, Brown-Saracino (2017) suggest that the spectrum of operational definitions "gesture to collective uncertainty about how to define and operationalize gentrification." The following section details the general framework that the majority of gentrification literature employs to identify the process. We highlight the limitations of these approaches, while also advancing a more holistic representation of gentrification within the broader context of neighborhood change.

Threshold Strategies

The majority of quantitative studies that identify gentrification have generally taken the form of a two-tiered, prescriptive approach in which spatial units are qualified as being eligible for gentrification during a starting year based on one or more socio-economic thresholds (figure 5). This eligibility step explicitly limits gentrification to lower socio-economic areas. As such, eligibility constraints typically require tracts to fall either below the regional median or the lowest quartile in one or more socio-economic measurements for consideration. Hammel and Wyly (1996) use an eligibility threshold requiring eligible gentrified census tracts to fall below the city-wide median income during a base year. Ellen and Ding (2016) require eligible tracts to be located in the central city with average family incomes below the 40th percentile of metropolitan-wide average family incomes. Maciag (2015) defined eligible tracts as those that were in the bottom 40th percentile of both median household income and median home value, compared to all tracts within the metropolitan area.

Areas that meet the eligibility precondition are then distinguished as gentrified by outpacing one or more proxy measurements for neighborhood change. Hammel and Wyly (1996) identify gentrified tracts as those that are eligible and rise to above city-wide median incomes in the subsequent decade. Ellen and Ding (2016) require a minimum 10 percentage point increase in the tract-to-metro ratio of average family income, percentage of white residents, percentage of college-educated residents, or median rent. Maciag (2015) requires gentrifying areas to experience both increases in college attainment and inflation-adjusted median home values that are in the top third percentile of all metro census tracts.

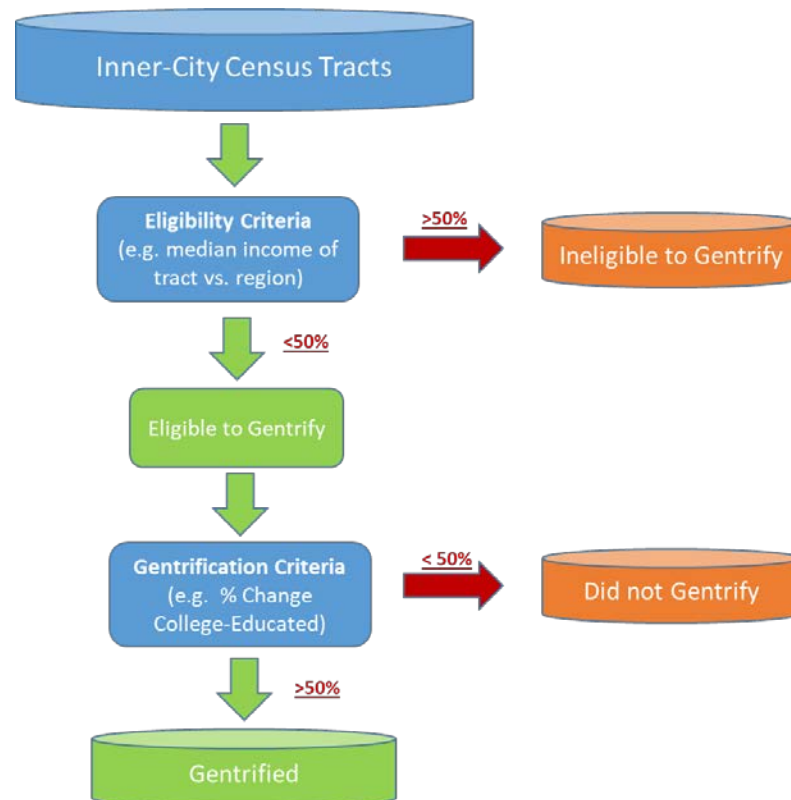


Figure 5. General prescriptive framework used to quantitatively identify gentrification.

Previous research built off of this model with varying input variables and thresholds, but maintained its general classification framework. The examples show the differences in arbitrary threshold values that can control the inclusivity or exclusivity with little or no explanation to justify their selection. Freeman (2005) adopts a median threshold approach with the disclaimer that “the median is an admittedly arbitrary threshold.” In addition to threshold values, the variables that are selected as indicators for gentrification vary considerably. Some studies consider demographic turnover using education, income, employment, and/or race. Investment patterns are frequently captured by different surrogate variables including home values, rent, mortgage lending patterns, and/or built environment renovations (Brown-Saracino, 2017; Hwang and Sampson, 2014; Kreager et al., 20011).

In addition to the inconsistent metrics applied within these frameworks, there are issues of confirmation bias and limits to the number of dimensions that they are capable of capturing. These prescriptive, threshold approaches suffer from confirmation bias by assuming *a priori* that gentrification is a significant neighborhood-shaping process within the region and time period of interest. The methodologies fulfill this expectation by comparing change measurements relative to city-wide medians or averages, which will almost invariably classify some areas as gentrifying regardless of the specified time period and city.

Threshold-based specifications cannot support the inclusion of multiple social, economic, and built environment variables. Adding additional criteria would essentially guarantee fewer areas identified as gentrified unless threshold values are broadened. This too often results in extremely crude measurements that stray from the actual complex process of gentrification. Furthermore, these methods do not attempt to account for inter-city differences of gentrification. Instead, they mistranslate an arbitrary and limited set of criteria to the multi-dimensional process of gentrification. The following section uncovers the red flags in the mainstream methods to identify gentrification. We conclude by pointing towards a promising new framework that bypasses the limitations of previous methods in order to distinguish gentrification apart from other types of neighborhood change.

Inter-definitional Disagreements

As early as 1986, Galster and Peacock questioned the sensitivity of differing operational definitions of gentrification that. They identified large discrepancies between the number and location of gentrified tracts based on differing criteria and thresholds. Their Philadelphia-based study identified 65 eligible census tracts with values less than four eligibility criteria: 1) median home value and 2) percent college educated below city-wide medians, 3) income less than 80% of city-wide median, and 4) percentage of white households less than 90% of the tract population. They then examined four simplified operational definitions used to identify gentrification between 1970 and 1980: proportion of black households, proportion of college educated, income change, and home value change. Each definition was evaluated against one of three stringency thresholds: low, medium, or high. The lowest stringency defined gentrification as any tract with a percent change greater than the city-wide mean change of that criterion during the decade. Medium stringency for each criterion was set to return half

the number of tracts of the low stringency value and high stringency was set to return only four gentrified tracts (Galster and Peacock, 1986).

The study showed that as many as 82 percent or as few as 6 percent of eligible tracts could be identified as gentrified based on changing the operational definitions and/or thresholds. Using only the lowest stringency (median-based change) 53, 17, 49, and 13 census tracts could be identified as gentrified based on the chosen operational definition alone. In addition, there was little correlation in the number of identified tracts by definition, with an average correlation of .29 for the medium stringency condition.

Enterprise Community Partners (2019) recently launched an online Gentrification Comparison Tool (GCT) to shed light on the methodological differences used to identify gentrification across 93 cities. The tool maps three gentrification studies' threshold-based definitions of gentrification detailed in table x.

There are important differences to highlight in each definitions approach to identify gentrification. Each definition uses a different threshold to classify areas that are potentially "gentrifiable" based on varying degrees of income. Freeman's (2005) definition uses the 50th percentile of median household income while Ellen and O'Regan (2010) use 70th percentile of average household income, both relative to metro-wide averages. The McKinnish et al. (2008) definition uses the 20th percentile of average family income compared to a national sample of urban tracts. Freeman's eligibility criteria also adds an additional prerequisite, using housing age as a proxy for disinvestment.

Table 1. Eligibility and gentrification criteria as applied by Freeman (2005), Ellen and O'Regan (2010), and McKinnish et al. (2010).

	Freeman	Ellen and O'Regan	McKinnish et al
Eligible	Median Household Income and percentage of housing built in prior 20 years both less than metro-wide values in 2000	Less than 70% of metropolitan average household income in 2000	Average family income in bottom 20% of nationwide urban tracts in 2000
Gentrified	Eligible and change in residents with college degree greater than metro-wide average and increase in real housing prices between 2000 and 2010	Eligible and minimum of 10 percentage point increase in the ratio of tract-to-metro average household income between 2000 and 2010	Eligible and real increase in average family income of at least \$10,000 between 2000 and 2010

The methodological differences continue with each definition's requirements for gentrification. The Ellen and O'Regan and McKinnish et al definitions both use different income measurements (household vs. family) to signal gentrification, as well as differing threshold types that are defined as relative for Ellen and O'Regan and absolute for McKinnish et al. Contrarily, Freeman only considers the cardinality of housing prices in addition to the relative gain in college educated residents as a proxy for socio-economic upgrading.

We used the GCT to evaluate the areas identified by these three competing definitions in Davidson County- Nashville, TN between 2000 and 2010. Figure 6 displays large disagreements between both the areas identified as eligible for gentrification and those ultimately classified as gentrified.

Freeman's definition immediately appears the most inclusive classification of the three approaches, labeling 37 census tracts as gentrified. Ellen and O'Regan's methodology identified 8 gentrified tracts while McKinnish et al identified 5 gentrified census tracts. Only three census tracts in Davidson County were consistently identified as gentrified across all three definitions.

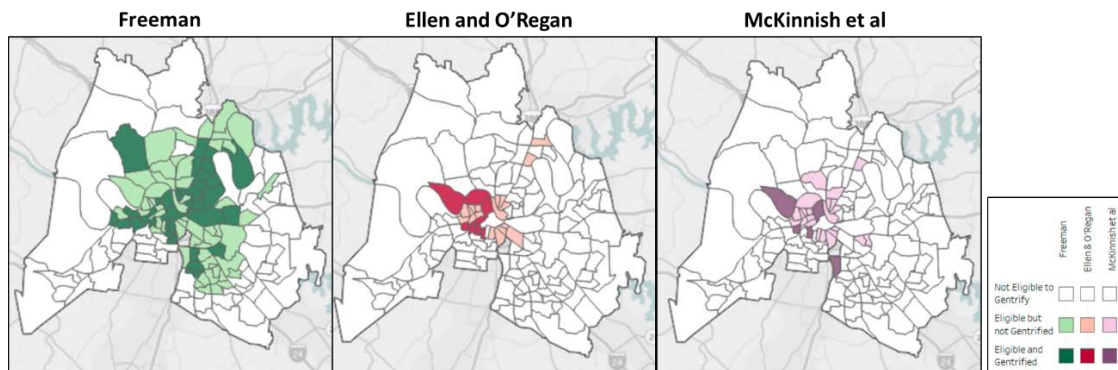


Figure 6. Variability in areas identified as gentrified within Davidson County, Tennessee (2000 - 2010).

The alarming inter-definitional disagreement becomes even more apparent when extrapolated across the 93 cities. Freeman's definition identified 3,221 census tracts as gentrified between 2000 and 2010 compared to 781 for Ellen and O'Regan and 501 for McKinnish et al. Only 152 tracts were consistently identified as gentrified out of 3,787 total tracts which were identified as gentrified by at least one method, amounting to a mere 4% agreement between definitions.

Intra-definition Sensitivity

The differences across threshold-based methods to identify gentrification are substantial, however there may even be significant variability within a single definition. This can occur within the specification of thresholds or cutoff values. We explore the sensitivity of Lance Freeman's definition of gentrification further in order to highlight how arbitrary threshold values may significantly affect the identification of gentrifying areas. Freeman's constructed definition requires gentrifying census tracts to fall below the metropolitan-wide median during the base year and also experience an increase in educational attainment above the metropolitan-wide median. Freeman explains, "the median is an admittedly arbitrary threshold." He reconciled this by reporting a second equally-arbitrary, but narrower threshold criteria, requiring eligible tracts to fall below the 40th percentile of the income and housing variables at the base year. Freeman's median-based definition identified 6.5% of total urban tracts as gentrifying, while the 40th percentile version identifies 3.9% of total neighborhoods as

gentrifying⁵. Extrapolated across his nationwide study area of 50 metropolitan areas and 45,108 tracts, this amounts to a difference of nearly 1,200 neighborhoods that are either included or excluded by a 10% difference in the selection of an arbitrary eligibility criteria.

We applied Freeman's competing 40th percentile and median-based threshold criteria for gentrification to Nashville between 1990 and 2000 in order to evaluate the potential sensitivity between the two thresholds. Figure 7 plots the results of the two methods. The median-based definition classified 56 tracts as "gentrifiable", of which 30 satisfied the definition's requirements for gentrification. The narrower, 40 percentile threshold identified 36 eligible tracts, of which 26 are labeled gentrified.

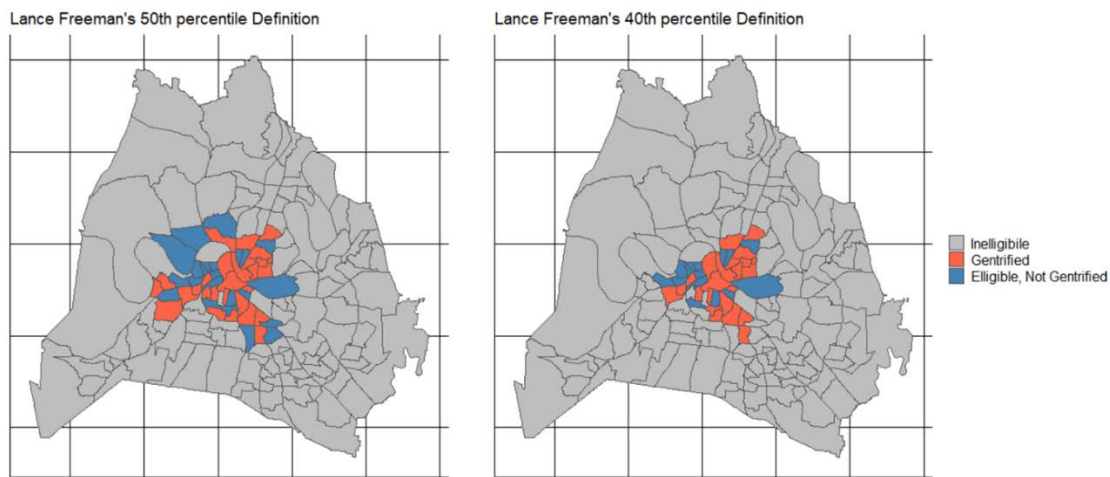


Figure 7. Variability between Freeman's (2005) 50 percent and 40 percent eligibility criteria for gentrification.

Here we explore where threshold-based definitions for gentrification may stray from the intended target in either direction by looking at two hypothetical examples applied to Freeman's median-based criteria for eligibility and gentrification. Consider an inner-city census tract that is eligible for gentrification because it falls into the 49th percentile of both household income and proportion of new housing in 1990. The tract outpaces the metropolitan percent increase in educational attainment by 1% and exhibits a real housing price increase of \$1. This tract would be considered stable by nearly all other accounts with the exception of the current constructed definition of gentrification. Another example considers a census tract in the 1st percentile of both of the eligibility criteria (lowest possible income and oldest housing), again satisfying the eligibility constraints. Real housing prices triple over the selected time period, but the percent increase in educational attainment falls one percentage point below the metropolitan median and the tract does not qualify as gentrified. It's in this context that we

⁵ Freeman's median-based definition also identified 41% of eligible tracts as gentrified between 1980 and 1990. Between 1990 and 2000 Freeman classified 31% of eligible tracts as gentrified using median-based criteria.

can plainly see that threshold based definitions may “doubtlessly commit errors of inclusion and omission” (Vigdor et al., 2002).

A replication study by Barton (2014) compared Freeman’s median-based methodology to the approach used by Bostic and Martin (2003). The latter technique identifies gentrification as those tracts with a median household income below the metropolitan statistical area (MSA) median in a base year (eligible) followed by income above the MSA median in the ending year as gentrified. Barton’s replication study identified a large discrepancy between both the number of tracts that are eligible and identified as gentrified. Between 1990 and 2000, 268 tracts were eligible for gentrification, of which 110 were classified as gentrified using Bostic and Martin’s (2003) definition of gentrification. Freeman’s (2005) definition yielded 582 gentrification-eligible tracts and 240 gentrified tracts.

Implications

We focus our attention on Freeman’s definition because of the real downstream implications of relying on an unstable definition of gentrification. A 2005 article by USA Today titled “Gentrification is a Boost to Everyone” summarized Lance Freeman’s findings, stating, “His conclusion: Gentrification drives comparatively few low-income residents from their homes” (Hampson, 2005). This conclusion rests on the assumption that gentrification was correctly identified upstream of his displacement analysis. As we have discussed and Barton’s (2014) replication study showed, this method to identify gentrification is highly sensitive to arbitrary threshold parameters. Barton’s results as well as the results from the GCT tool both suggest that Freeman’s methodology is a highly inclusive identifier of gentrification. This overly-inclusive definition of gentrification may pull in higher socio-economic tracts (up to the median household income) that could offset higher displacement rates in other tracts. Newman and Wyly (2006) re-examined New York City displacement data and concluded that 6-10% of all rental moves were affected by some form of displacement. They highlighted the discrepancy of Freeman’s control (non-gentrified) group that included some of the lowest income areas, effectively creating an artificially high comparison of displacement rates between gentrifying neighborhoods, fitting with the second hypothetical example outlined above.

Freeman’s operational definition of gentrification has been adopted for a wide range of studies in gentrification literature. Lester and Hartley (2013) use the Freeman’s definition to herald the local employment benefits in gentrifying neighborhoods. It is often used as a pillar of evidence against government interventions to curb the gentrification process. The threshold-based strategy to identify gentrification is not unique to Freeman, but is borrowed across the academic literature with slightly different proxies and thresholds. Because no one standard exists in the real world to qualify an area as “gentrifiable”, threshold values can too easily act as arbitrary levers to identify gentrification for an end-purpose.

Clustering Strategies

The limitations of previous quantitative methods to identify gentrification have spawned a more recent class of models to identify the elusive process. K-means clustering is a type of unsupervised machine learning technique used to distinguish underlying patterns and similarities among multiple features. In the context of urban change, it can be used to distinguish the dominant trajectories of neighborhood change over time. This method is compatible with Brown and Saracino's (2017) call "to study the city more holistically, capturing neighborhood change and stasis, poverty, affluence, and everything in between."

Clustering techniques bypass many of the limitations and bias of traditional threshold-based methods. Clustering methods optimize the homogenous subgrouping from the data itself, rather than from a predetermined framework. They also afford the inclusion of additional variables that are able to capture the multi-dimensionality of neighborhood change. These methods are also less biased than their prescriptive counterparts, because they do not make a priori assumptions about the underlying phenomena. Traditional methods, by comparison to city-wide values, have ensured that gentrification would be measured irrespective of the city and time period analyzed. The clustering method avoids this confirmation bias; it makes no guarantees that a cluster resembling gentrification will emerge, nor any other pre-determined type of neighborhood change for that matter. Rather the approach reframes the question from "Where is gentrification happening?" to "What are the emergent trajectories that neighborhoods change over time?"

This alternative methodology represents a fundamental shift in the burden of translating such a complex and polarizing process as gentrification to a set of pre-determined criteria. Whereas previous literature defined up-front thresholds to identify gentrification, the machine-based classification method distinguishes the major pathways in which neighborhoods change, and subsequently necessitates an interpretation of the observed trajectories. Using domain knowledge of the trends associated with gentrification, one can examine cluster outcomes to identify neighborhood change typologies and determine if any one cluster exhibits change characteristics associated with gentrification. The determination of a gentrification typology can then be defended by data-driven statistics, spatial organizations, and the wealth of domain knowledge advanced by previous researchers specific to gentrification. These changes are marked by increases in property values (both home and rent values) as well as demographic and socio-economic shifts reflected in income, education, and race. The output is a robust classification that may be used as a dependent variable for downstream analyses and prediction.

These advanced clustering techniques have only fairly recently been applied to the context of urban change. Podagrosi et al. (2011) used k-means clustering to distinguish neighborhood change between 1980 and 2000 in Houston, Texas. The researchers used principle component analysis to reduce a set of thirty-eight change variables to five dimensions of change. They identified fifty-four census tracts exhibiting similar levels of upgrading in line with gentrification. Across these tracts college graduates increased by 88%, home values increased by 30%, and per capita incomes increased by 53.1%.

Ling and Delmelle (2016) preprocessed eleven census change variables (1970 -2010) using a self-organizing map procedure before conducting a k-means clustering analysis to analyze temporal

trajectories in eight US cities. Ten different neighborhood typologies were identified: blue-collar suburbs, aging middle-class suburbs, struggling urban, suburban decline, struggling older suburban, early revitalization, decline then recovery, stable elite, suburban densification, and suburbanization. The neighborhood types also exhibited significant spatial autocorrelation across each city when measured using a join-count statistic.

K-means clustering has also been used to distinguish discrete neighborhood typologies as opposed to change typologies in Chicago and Los Angeles. Census tracts were classified as either elite, blue collar newer suburban, older stable suburban, struggling, or young urban for each decade between 1970 and 2010. The sequences of each tract typology were then analyzed according to a sequential pattern mining algorithm. The transition sequences fell into either stable, downgrading, or upgrading processes. Upgrading was found in central-city neighborhoods of Chicago, but was limited to the suburbs in Los Angeles (Delmelle, 2016).

We employ k-means clustering on unlabeled data to observe the natural neighborhood transitions that are inherent to Nashville and argue that one specific typology is reflective of an underlying gentrification process. The classification itself may provide valuable insight into the city-specific nature of neighborhood change. This method allows for a more comprehensive feature space on the front end of analysis to determine areas potentially undergoing gentrification. Additionally, this dependent variable can be easily deconstructed to find city-specific patterns that coalesce to form dominant the dominant pathways that neighborhoods. This methodology of defining gentrification can be used for predictive purposes, but assumes that past gentrification is indicative of the future process.

Nashville

The 2018 Nashville metro area population totaled over 670,000 people, representing the 24th largest incorporated city in the United States. Despite a fast population growth and anecdotal accounts of gentrification, Nashville has been omitted from a specific analysis of neighborhood change (Haruch, 2014; Larsson, 2017; Plazas, 2017).

Nashville claims the moniker “Music City” and is considered to be the country music capital of the world. Over the past 60 years the music industry expanded considerably beyond the country genre and stimulates substantial economic growth in Nashville. Nashville’s booming music industry supports more than 56,000 jobs and contributes \$5.5 billion to the local economy (Harper and Cotton, 2015). The success of the music industry has undoubtedly spilled over to the tourism sector which saw 15.2 million visitors to Nashville in 2018 (*Nashville Tourism & Hospitality*, 2018). These highly visible industries have gained Nashville a reputation as a leading cultural and creative center that attracts top talent in other industries. Several academic institutions including Vanderbilt University, Belmont University, Lipscomb University, Fisk University, Tennessee State University, and Meharry Medical College help support the pool of skilled labor comprising the Nashville economy.

Nashville has quickly become a destination city for transplants and companies alike. Big-name companies including Amazon recently announced plans to build an operations center in Nashville. Financial companies are increasingly eyeing Nashville, citing the low cost of real estate, taxes, labor, and utilities. Nashville ranked second in a study of 40 US cities financial industry (McGee, 2018). Nashville also hosts a thriving healthcare industry which contributes \$46.7 billion to the local (*Nashville Region's 2018 Vital Signs*, 2018).

Several high-level economic indicators suggest that Nashville's growth-oriented approach has been successful. The Nashville metropolitan area grew by 2.24% between 2000 and 2010, representing the seventh fastest growing city in the United States over this same time (Forbes, 2018).⁶ The Nashville metropolitan area grew by an average of 5,500 people per year between 2000 and 2010 before ramping up to 10,700 people per year since 2010 (*Housing Nashville: Nashville & Davidson County's Housing Report*, 2017). This population growth was supported by an influx of roughly 11,400 jobs per year between 2002 and 2014. Nashville wages also saw increases from \$35,000 per year to \$53,000 per year between 2000 and 2015. Additionally, the Nashville housing market appears to have fervently rebounded from the 2007 recession. Nashville also topped Zillow's list of the hottest housing markets in 2017 (Allison, 2017).

Costs of Growth

A closer examination of Nashville's rapid growth realizes troubling signs of inequality, unaffordability, and concerns over gentrification-induced displacement. Although wages increased annually by 2.8%, the Consumer Price Index cost of living metric for urban southerners also increased by 2.3% annually over the same period. This offset results in the average workers wage increasing only a modest .5% since 2000. Additionally, median household income increased by 1.3% per year between 2000 and 2014, but was outpaced by inflation (2.3%) over this time period. This resulted in an inflation-adjusted 1% per year decline in median household incomes.

Nashville saw more new housing constructed in 2015 than at the height of the housing bubble in 2007, however, this increase in new housing supply has not translated to more affordable housing options in Nashville (Fraser, 2017). Although some states encourage or require localities to increase the supply of affordable housing, others have preempted local efforts. In Tennessee, a recent state law effectively undermined the Nashville Metropolitan Council's inclusionary ordinance by requiring that the city provide financial incentives to developers that voluntarily include affordable units.

Both production and consumption side forces have impacted the affordability of housing in Nashville. A recent trend suggests that out-of-state investment firms have targeted Nashville and its growing suburbs, owning at least 4,900 homes in the region. Often times these homes are taken off the market for several years, reducing supply and creating artificially higher prices in an already unaffordable market (Reicher, 2017). Other forces acting on the consumption side of Nashville's

⁶ Ranking is out of the 100 largest U.S. metropolitan statistical areas.

affordability crisis are the effects of short term rentals like Airbnb, VRBO, and HomeAway on constraining the supply of housing. There are over 5,000 active short term rental permits in Davidson County, 2,700 of which are operated under Airbnb (*Airbnb, Housing, and Nashville*, 2017; Metro Codes Department, 2019). The growth in short term rental permits is provided in figure 8.

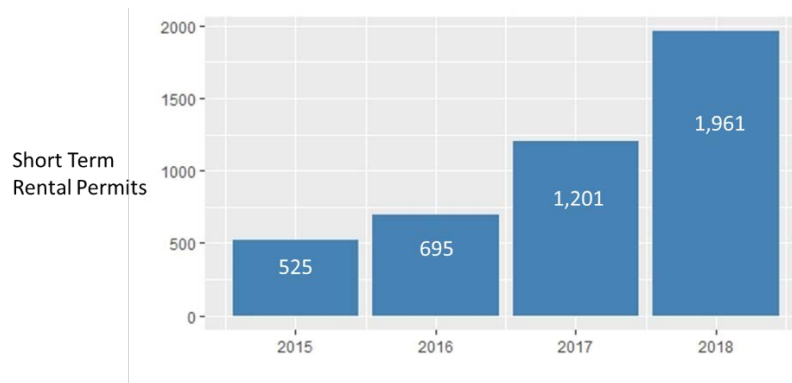


Figure 8. Davidson County short term rental permits (2015 - 2018) (Metro Codes Department, 2019).

These community concerns were highlighted by David Plazas (2017) of the Tennessean newspaper in a series of columns titled “Costs of Growth and Change in Nashville.” Many of Nashville’s longtime citizens report facing financial pressures amid rising rent and property tax values as well as social isolation as neighborhoods are redeveloped to a more expensive clientele. Plaza’s captivates the inequality and community concerns amid Nashville’s growth, writing “African Americans have been hit exponentially hard in Nashville.” He also cites the changing socio-economic geography amid gentrification: “The segregated areas where [African Americans] once lived around the urban core, like East Nashville, Germantown, and Edgehill, have now become high-rent, whiter communities” (Plazas, 2017).

Data Sources

The United States Census is one of the most important data sources within the social sciences. The time period of our analysis necessitates the cross-comparison of US decennial census (DC) data for 2000 data and the American Community Survey (ACS) for 2016 data. Between 1960 and 2010 the long form questionnaire of the DC was distributed every decade to approximately one in six US households with the remaining households completing the short form DC. Post-2010 the ACS supplemented the long form DC with the consistent goal of collecting detailed demographic, economic, and housing data.

There are fundamental differences between the sample size and timing of the long form DC (2000) and ACS (2010- present) data sources that should be addressed. The 2000 long form DC was distributed to 18 million households whereas the ACS comprises roughly 3 million households annually. Our analysis uses the 2012 – 2016 ACS 5-year estimates and should be understood to represent the rolling average over a 5 year time period as opposed to the DC survey which is representative of a single year.

The small sample size of the ACS has the potential to diverge from the ground truthed measurements. It is possible that what the ACS gains in temporal granularity may be countered by a loss in accuracy. Bazuin and Fraser (2013) resampled a single census tract in Nashville, TN using the DC approach and found a gross underestimation of both total population as well as the number of people living in poverty.

Another cautionary difference between the two surveys is the reference period for employment and income questions. The ACS survey asks employment status relative to the week before the survey completion while the DC uses a static employment reference period of the week before Census Day (April 1, 2000). Similarly, ACS income and earnings data is reflective of the previous 12 months, while the DC asks income in the previous calendar year. These methodological differences are the result of the year-round surveying of the ACS versus the two month window between March and April when the DC survey results are returned (US Census). An additional criticism of the ACS is the potential underrepresentation by disadvantaged groups.

Census data from the year 2000 were obtained from the Neighborhood Change Database (NCDB) (GeoLytics, Inc., 2005). The NCDB apportions U.S. decennial census data to 2010 spatial boundaries. 2016 ACS data were collected from the American Community Survey API via the R package *tidycensus* (U.S. Census Bureau, 2016; Walker, 2019). Spatial geometries were accessed from the *tigris* R package (U. S. Census Bureau, 2017; Walker and Rudis, 2019).

In alignment with reproducible research principles, all data cleaning, restructuring, and statistical analyses were conducted in R with the supplemental code submitted in the final work product. Several open-source R libraries were critical for carrying out the analyses in this paper including *stats*, *caret*, *randomforest* and *dplyr*, among others (Kuhn, 2019; Liaw and Wiener, 2018; Wickham et al., 2019).

A GIS shapefile containing Davidson County park boundaries was obtained from the data.nashville.gov data portal (Metro Government of Nashville & Davidson County - Parks and Recreation, 2016). The ArcGIS tabulate intersection tool was run on the park polygons and 2010 Census tract boundaries in order to calculate the percentage of census tract occupied by park space.

Other data sources include building permits issued by local government. A permit database is provided by nashville.data.gov through an open data portal (Metro Codes Department, 2018). This dataset contains building permits dating from 2013 to present and is updated daily. The building permit database provides the location of the permit, permit type (Demolition, Residential-New, Commercial-New, Residential-Rehab, Commercial-Rehab, etc.), and the construction cost associated with the permit. This detailed dataset can be used to gain insights into capital investment types believed to accompany patterns of gentrification and neighborhood change. Additionally, eviction data was used as a predictive variable from Eviction Lab (Desmond et al., 2018).

Scope

This research focuses on Davidson County, Tennessee. The urban center of Nashville is located approximately in the geometric centroid of the county. Census Tracts were used as the primary unit of analysis for reporting. 161 census tracts comprise the study area of Davidson County, Tennessee. Census tracts are nested within county boundaries with populations of approximately 4,000 residents. Tracts were delineated by the U.S. Census Bureau to provide relatively stable boundaries used in the statistical analysis and comparison of Census data across time (U.S. Census Bureau, 2018). City planners regularly appropriate “neighborhood” boundaries that are roughly in line with Census tract boundaries (U.S. Census Bureau, 2018). Census tracts were also the spatial unit chosen for the Nashville Metro Council’s Housing Policy and Inclusionary Zoning Feasibility Study (2017).

The selected time period for this study is motivated by figure 9, indicating the average residential property transaction distance to downtown since the 1950’s. The start of the 21st century coincides with a major turning point between suburbanization and a back to the city movement. The year 2000 also coincides with the release of the decennial census report and will serve as the base year for this analysis on neighborhood change. 2016 marks the most recent release of American Community Survey data and is used as the end year by which trends are analyzed. The 16 year time period is also in line with previous gentrification literature such as Goetz et al. (2019) which focused on identifying the gentrification in Minneapolis between 2000 and 2016. The selected time period is also less sensitive to the housing crisis and economic recession between 2007 and 2009.

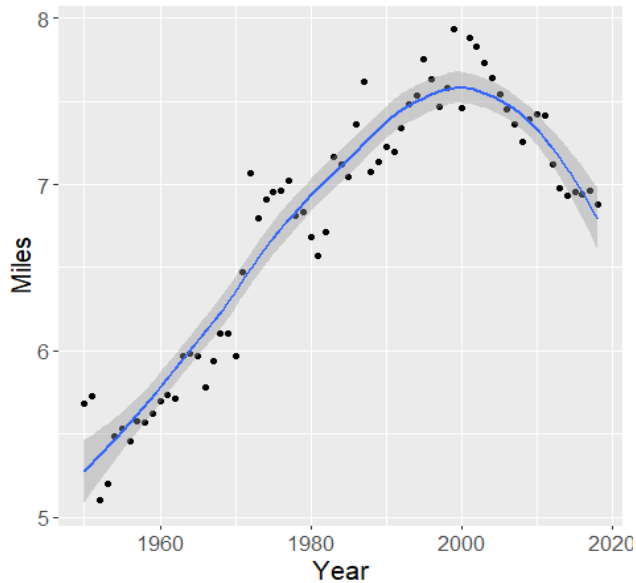


Figure 9. Davidson County home sales as a function of distance to downtown (Davidson County Tax Assessor, 2018).

Identifying Gentrification Through Machine Learning

Machine learning (ML) is a broad class of powerful statistical techniques that relies on artificial intelligence (AI) to learn from past observations to perform a specified task without being explicitly programmed. ML has made substantial advances in fields such as biomedical, image processing and speech recognition due to the vast complexity and relationships within real-world problems. To a smaller extent, the same holds for gentrification and the processes that lead to neighborhood change. The majority of quantitative studies identifying and predicting manifestations of gentrification rely on statistical modelling techniques that include logistic modeling and linear regression (sources) to predict a user-defined dependent variable of gentrification. Several barriers precluded early urban researchers from implementing ML techniques that include computational limitations as well as a lack of available data and gentrification outcomes.

An agreed-upon, quantitative definition of gentrification does not exist. We previously outlined the limitations, sensitivity, and biases that follow quantitative attempts to identify gentrification. We present a methodology using six census variables (rent value, home value, income, race, educational attainment, and percentage of multi-unit buildings) commonly used to identify gentrification. The six variables are calculated as percent change between 2000 and 2016. By looking exclusively at percent change, we are interested in the natural evolution of neighborhood change, regardless of their starting socio-economic or demographic compositions. Home and rent values are perhaps the best metrics

capture investment and disinvestment. Because gentrification has different consequences for renters versus homeowners, both rent and home values were included. Built environment changes are also signaled by changes in the percentage of multi-unit dwellings. Income, race, and educational attainment changes are included to capture the flows of different classes of residents that could potentially signal displacement.

Previous researchers have used larger number of variables to identify neighborhood change trajectories (Delmelle, 2016; Podagrosi et al., 2011). To accommodate these variables, they often use principle components analysis to reduce the dimensionality down to a handful of proxy variables that attempt to measure the same fundamental aspects of the built environment changes. We cluster on only six of the most important proxies for investment, disinvestment, and demographic change. The six features selected for clustering strikes a balance of more dimensions than the prescriptive methods, but less features than previous clustering methods. This bypasses the need for dimensionality reductions like principle components analysis and allows for more direct and intuitive reporting without sacrificing or diluting the most important factors.

Perhaps the most time and labor intensive phase of any ML analysis involves gathering, cleaning, aggregating, and selecting relevant data. We have compiled a dataset of 16 variables that we believe may provide insight into trajectories of neighborhood change. As such, we employ several ML techniques both supervised and unsupervised to mine the dataset for patterns and predictions that may benefit public understanding, decision-making, and future resource allocation.

We err on the side of interpretation by selecting only six census-based change variables to identify a gentrification typology. Additional variables could be added to the front-end of this analysis, which may provide additional insight into the underlying structure of neighborhoods change. At a certain threshold it may be required to preprocess Principle Component Analysis to preprocess and reduce the dimensionality of the identification variables into the chosen clustering technique. This theoretically sacrifices interpretability of the outcomes, but affords the considerations of more variables indicative of gentrification.

K-Means Clustering Neighborhood Change

K-means clustering is a type of unsupervised clustering algorithm that partitions observations into K number of user-specified groupings. The k-means objective function iteratively assigns observations to a cluster that satisfies the minimum within-cluster sum of squares (MacQueen, 1967). The proceeding step calculates and adjusts the new means to the centroids of the new cluster. We performed a K-means clustering analysis on the percent change between 2000 and 2016 of six census variables listed in table 2.

Table 2. Six socio-economic proxy variables used in the k-means clustering procedure (calculated as percent change 2000 - 2016).

Clustering Variables
Median Home Value
Median Rent Value
Median Household Income
Percentage of persons 25+ with college degree
Percent Non-White
Percent Multi-unit housing

K-means can be sensitive to the starting location of cluster centroids (Yi et al., 2010). We initialized the algorithm by selecting the optimal starting centroids out of 50 random permutations and report little variability within the initialization selection.

The selection of the optimal number of clusters in the dataset is a fundamental problem in cluster analysis (Sugar and James, 2003) and specific to the problem context. As such, we report direct methods of total within-cluster sum of squares (WSS) and average silhouette scores to guide our selection of the optimal number of clusters. WSS can be thought of as a measure of compactness, with the objective that the selected number of clusters minimizes intra-cluster variation. Plotting WSS against the number of clusters (figure 10) shows the diminishing return of WSS improvement as additional clusters are added. The elbow method is used to visually ascribe the optimal number of clusters that occurs when the WSS improvement from adding additional clusters levels off. Figure 10 shows a small drop-off in WSS accuracy between k of 4 and 5. Additionally we validate the choice of k via the silhouette method, which compares a record's similarity (Euclidian distance) to its own cluster against other clusters. K= 2 maximized average silhouette widths, but would suggest a binary classification of neighborhood change that does not fit the objective of this research. K =4 was chosen as it proved the second highest silhouette width as well as a subtle plateau in the elbow method. The selection of 4 clusters is also in line with Binet (2016) who identified 4 groupings as the optimal number of pathways of neighborhood change in Queens County and New York County, New York.

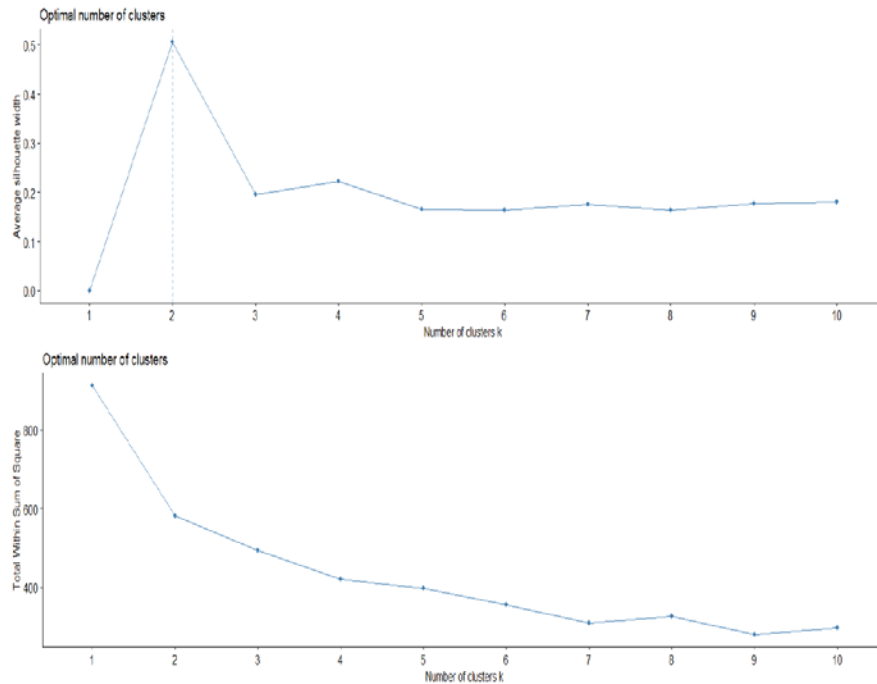


Figure 10. K-means clustering K-selection metrics (silhouette width and total within sum of squares).

Figure 11 visualizes the cluster geometries associated with different k values plotted against the first two principle components that explain 74% of the variance in the dependent variable dataset.

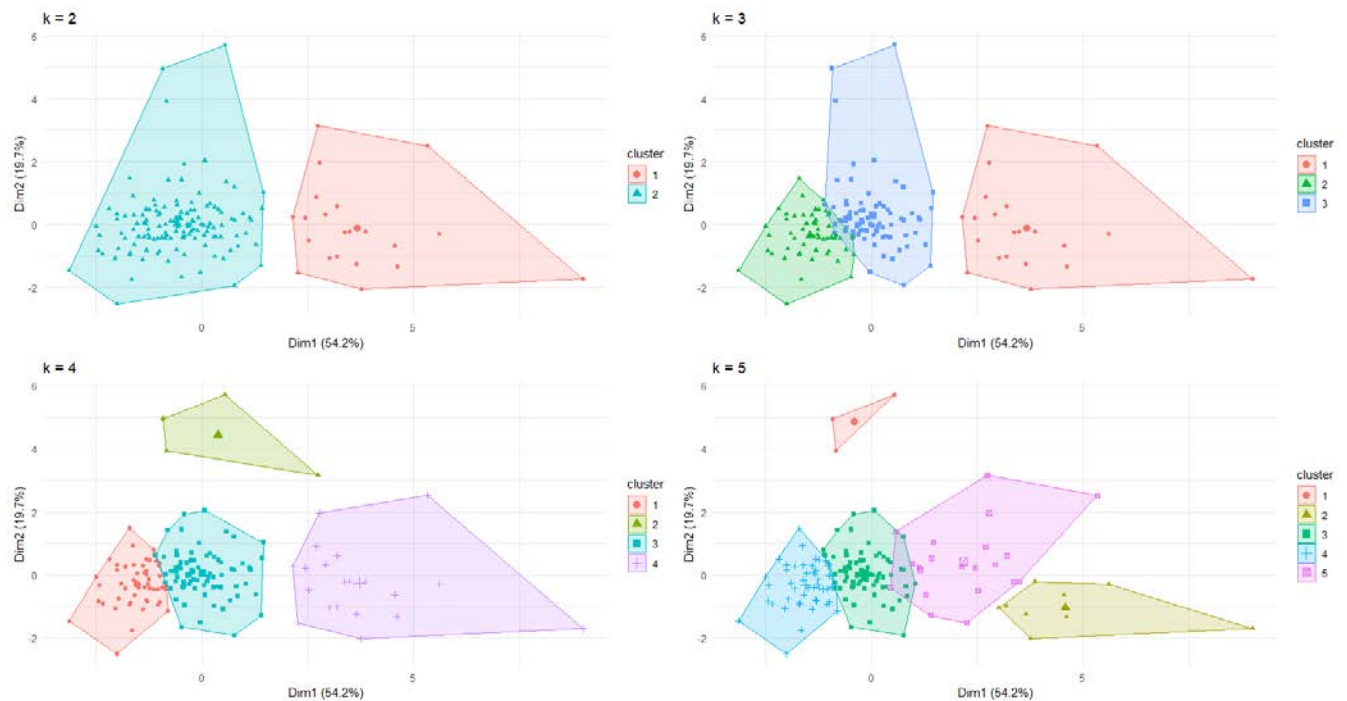


Figure 11. Cluster geometries based on cluster size specifications. Plotted against the first two principle components (73.9% of total dataset variance).

Nashville Neighborhood Change

We investigated neighborhood change using k-means clustering of percent changes across six census-based change variables between 2000 and 2016. Fundamental to this procedure is the selection of the appropriate number of clusters. We selected a cluster size (K) of four to balance both classification and accuracy. For the purpose of this study we were interested in identifying emergent typologies that may be indicative of gentrification. The socio-economic characteristics of gentrification are distinct from other patterns of neighborhood change; they are marked by increases in housing prices (home and rent values), socio-economic changes (income, education, race), and shifts in the built environment (multi-unit housing). Therefore, a carefully interpretation of the results of the k-means procedure is necessary. We build our interpretation of the cluster results from spatial distribution of clusters, summary statistics of the raw values and the domain knowledge of neighborhood change characteristics advanced by previous researchers.

Table 3 provides the raw summary statistics for county-wide change as well as the four emergent neighborhood change typologies. Before consideration of individual clusters, it is valuable to evaluate the broader patterns of change of our selected variables across all census tracts in Davidson County. County-wide median home and rent values increased by 21 and 13 percent, respectively. Inflation-adjusted household income saw an average 5 percent decrease. The proportion of multi-unit housing fell by a modest 1.1 percent average. Lastly, both educated and non-white percentages grew by 7.4 percent and 4.8 percent across Davidson County. These measurements provide the baseline and context to evaluate neighborhood change typologies.

Table 3. Cluster Summary Statistics (K=4).

	Mean	Median	Min.	Max.	SD
Davidson County (n = 153)					
Home Value	21.2	7.0	-38.4	246.8	43.7
Rent Value	13.4	7.6	-42.0	145.7	26.5
Household Income	-5.0	-10.9	-46.3	145.2	28.9
Percent College	7.4	5.2	-14.4	48.5	10.4
Percent Non-white	4.8	5.2	-54.7	31.1	13.6
Percent Multi-unit Housing	-1.1	-0.5	-68.4	37.4	12.4
K1 (n = 50)					
Home Value	-4.7	-5.8	-38.4	34.7	14.8
Rent Value	-6.2	-4.0	-42.0	14.7	12.0
Household Income	-22.8	-24.2	-46.3	11.9	11.0
Percent College	0.4	-1.2	-14.4	14.8	6.3
Percent Non-white	14.9	15.0	-8.2	31.1	9.5
Percent Multi-unit Housing	4.9	4.2	-13.7	37.4	9.3
K2 (n = 20)					
Home Value	108.4	106.5	39.1	246.8	54.3
Rent Value	45.7	38.8	6.6	145.7	30.3
Household Income	47.1	27.5	-11.9	145.2	43.3
Percent College	25.4	24.8	13.7	48.5	9.1
Percent Non-white	-15.6	-18.4	-54.7	6.6	17.0
Percent Multi-unit Housing	-4.0	-3.6	-18.6	8.9	8.0
K3 (n = 79)					
Home Value	16.7	13.3	-15.5	73.7	20.5
Rent Value	15.6	10.5	-11.3	103.7	20.9
Household Income	-7.6	-9.2	-46.1	29.1	13.0
Percent College	7.7	7.0	-5.3	27.0	6.5
Percent Non-white	3.0	2.7	-12.2	28.0	6.6
Percent Multi-unit Housing	-1.5	-1.6	-28.3	18.5	8.2
K4 (n = 4)					
Home Value	-3.2	-12.7	-18.7	31.2	23.2
Rent Value	53.3	46.1	34.4	86.6	23.6
Household Income	9.6	8.0	-6.3	28.6	18.1
Percent College	0.1	-1.7	-13.2	16.9	12.8
Percent Non-white	14.7	18.9	-3.9	24.9	12.8
Percent Multi-unit Housing	-51.4	-50.3	-68.4	-36.6	14.5

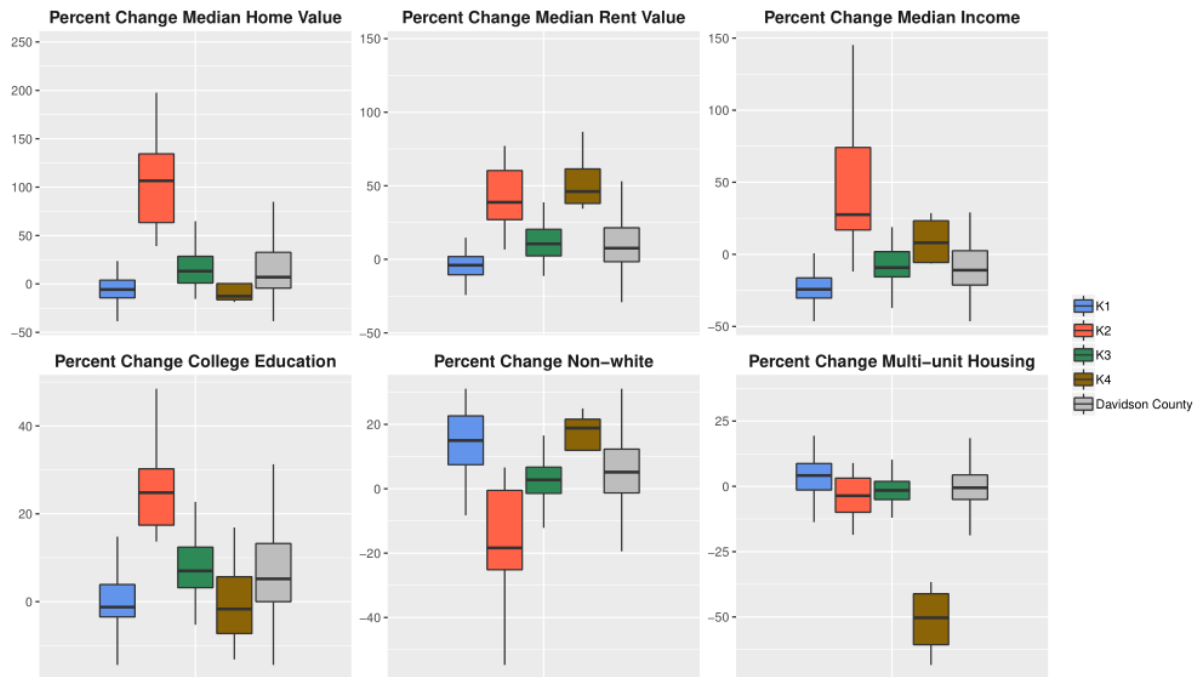


Figure 12. Box and whisker plots of neighborhood change typologies.

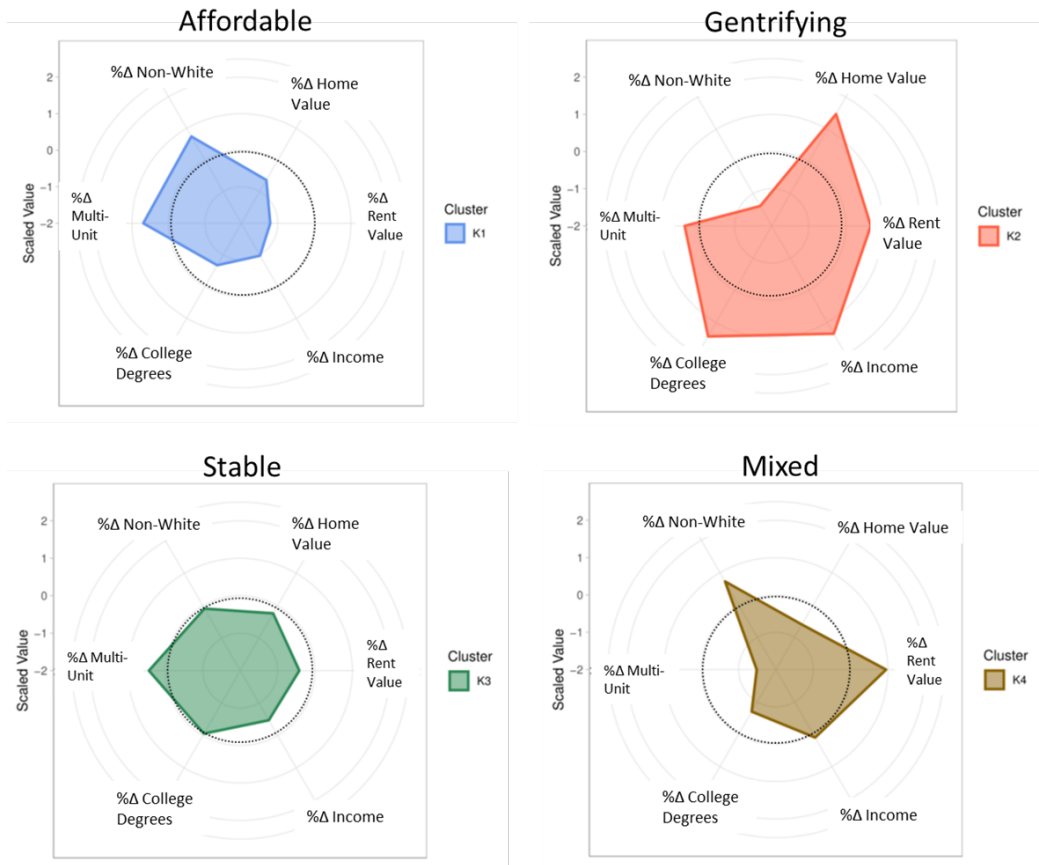


Figure 13. Davidson County normalized neighborhood change typologies ($K=4$).

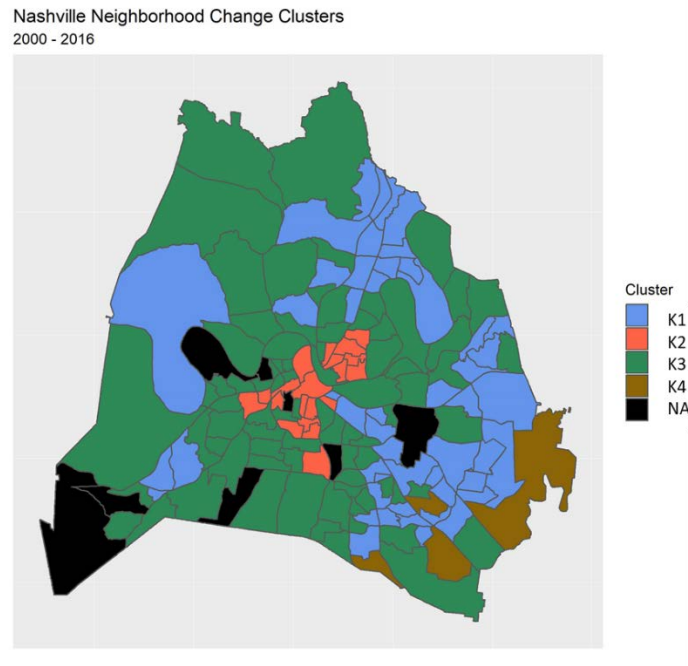


Figure 14. Davidson County neighborhood change cluster map (K=4)

The first cluster, K1, comprises 50 census tracts. We argue that this neighborhood change typology is indicative of disinvestment, racial diversity, and a possible beacon of housing affordability. This cluster reports the largest absolute decreases in home (-4.7%), rent (-6.2%), and income (-22.8%) measurements, all well below Davidson County averages. Educational attainment increased by a modest .4 percent, but lags behind the Davidson County average. The percent change of multi-unit housing within this cluster is a positive 4.9 percent increase, the highest increase in any cluster. This neighborhood change typology is also distinguished by a 14.9 percent increase in non-white population, representing the largest increase in racial diversity across any cluster. This neighborhood change typology is spatially distinguished in northeastern and southeastern corridors of the county.

We interpret cluster 2 as a polarizing cluster that exhibits the distinguishable patterns of gentrification and evidence of displacement. This typology is comprised of 20 census tracts spatially clustered near downtown Nashville. It is differentiated by the largest increase in home value (+108%), education (+25%), and income (+47%) across all tracts. This typology simultaneously exhibits the largest decrease in non-white occupants (-16%) across all clusters, breaking with the county-wide 4 percent increase in non-white individuals. The average percent decrease in multi-unit dwellings (-4%) falls below the Davidson County average decrease (-1%) indicating a restructuring within the built environment and housing composition. Cluster 2 is situated compactly around downtown Nashville, meeting the central-city criteria that is often reported in previous efforts to distinguish gentrification. Figure 15 reports the precise locations of the gentrification typology of neighborhood change (K2). Four reaches of gentrification are observed: East, West, Central, and South (relative to downtown Nashville).

Cluster 3 encompasses 79 census tracts and is interpreted as relatively stable tracts. This typology comprises the largest cluster of any neighborhood typology and falls consistently closest to the

Davidson County median values across all six input variables. K3 is spatially the most expansive cluster, encircling the gentrification typology and extending north, south, and northeast to the edges of Davidson County.

The final cluster, cluster 4, contains only 4 census tracts. It is differentiated by high rent increases, but overall decreases in home value. It also comprises the largest decrease in percentage of multi-unit housing. K4 appears as smallest spatial extent and limited to the Southeastern portion of the county.

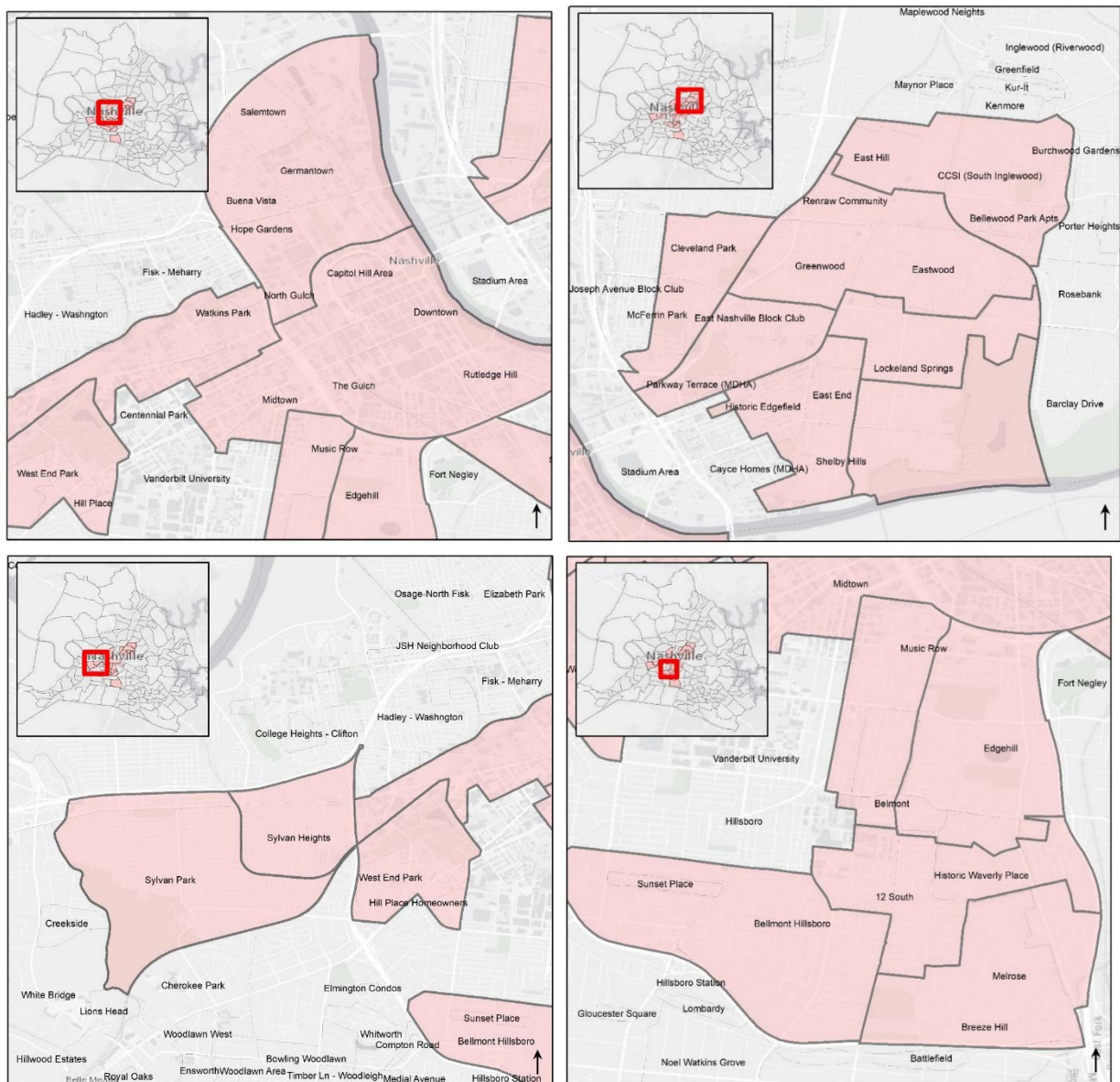


Figure 15. *Zoomed locations and neighborhoods of gentrification typology.*

Sensitivity

Figure 16 shows the mapped sensitivity regarding the selection of K, number of clusters. A distinct gentrification typology is apparent in K=2, K=3, and K=4. The occurrence of a gentrification typology at only k=2 further suggests that gentrification is one of the most dominant and distinguishable patterns of neighborhood change in Nashville between 2000 to 2016. Interestingly, at k=5, a portion of the tracts that are stably identified as gentrified (between K=2 and K=4) are grouped into a fifth cluster. This additional cluster contains tracts with more moderate property value increases as well as lower decreases in non-white populations compared with the remaining gentrifying areas (k2). This suggests that the selection of K clusters may be relatively insensitive up to four clusters, after which a portion of the gentrified areas are more closely related to other areas and their associated change pathways. More broadly, the specification of a larger K could allow for differentiation between different types of gentrification (moderate/advanced, early/mature, or displacement/incumbent upgrading).

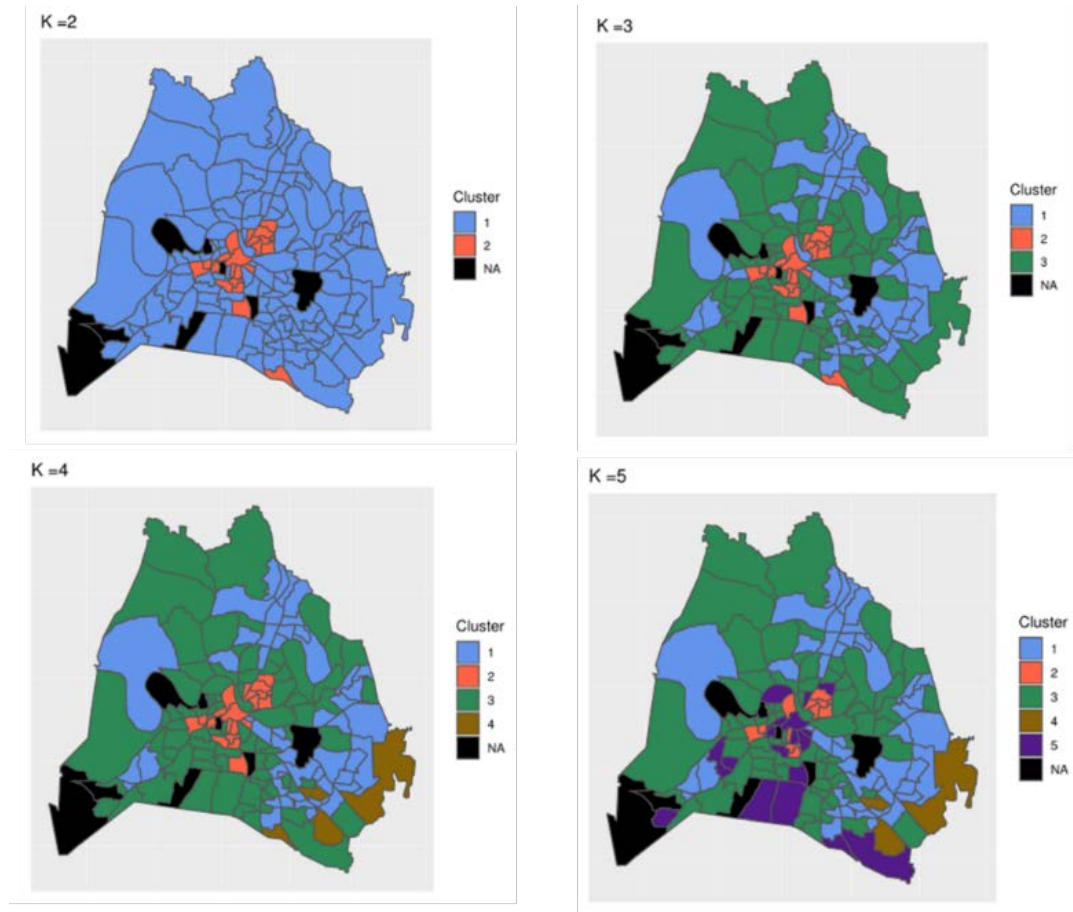


Figure 16. Davidson County K-means clustering (K) sensitivity (K2:K5).

K-Color Join Count Statistic

Spatial autocorrelation is a measurement of the spatial relationship of a variable with respect to its spatially-defined neighbors. This concept builds off of Tobler's first law of geography (1970), "Everything is related to everything else, but near things are more related than distant things." Spatial autocorrelation may take three forms: positive, negative, or zero. Positive spatial autocorrelation describes a statistically significant spatial clustering of similar values, while negative spatial

autocorrelation describes the spatial clustering of dissimilar values. Zero spatial autocorrelation implies that the observed spatial ordering is not statistically significant or distinguishable from random chance (Cliff and Ord, 1973).

We test the spatial independence of the categorical k-means clustering outcomes using the k-color join count statistic developed by Dacey (1968). This statistical test considers a null hypothesis that the observed local distribution of observations are random. 10,000 Monte Carlo permutations were run using the “joincount.mc” function in the “spdep” R package (Bivand, 2019).

We formally define neighbor relationships through two separate adjacency matrices: first order queen’s case and second order queen’s case. The first-order queen criterion defines spatial neighbors that share either a common edge or common vertex while second-order extends the neighbor relationship by including common neighbors that also share an edge or vertex. Queen’s case is the preferred contiguity for irregularly shaped aerial units like census tracts (Anselin, 2018). Figure 17 shows the resultant adjacency matrices between the first and second-order queen contiguity. The first-order contiguity case is a more localized measurement of adjacent neighbors with an average of six links per census tract. Second order contiguity identifies 18 neighbors on average per census tract. For sensitivity, we also examined rook’s case contiguity (requiring a shared edge); First-order rook’s case contiguity resulted in an average 5.07 links and did not affect our overall interpretation of spatial autocorrelation.

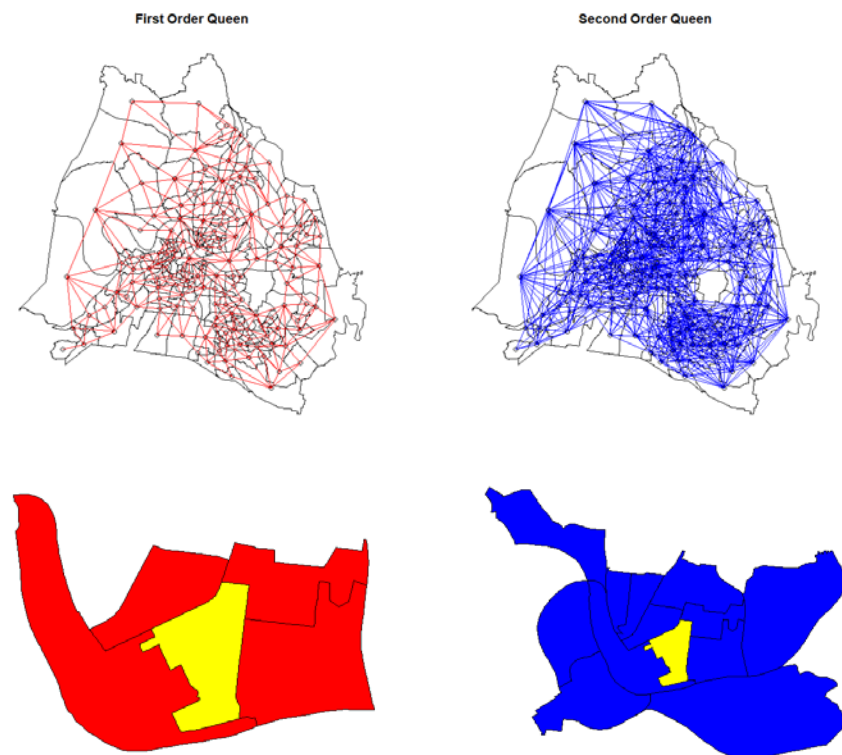


Figure 17. Comparison of first and second-order queen's adjacency matrix and sample census tract.

We evaluated the spatial autocorrelation of the categorical distribution of clusters using the join count test in order to test a null hypothesis that the observed spatial distribution was due to chance in figure 14. The k-color join count test returns p-values for each nominal variable's spatial distribution (table 4); First-order queen's case contiguity resulted in p-value < .05 for all clusters with the exception of K3 which appears exclusively in the Southeastern corner of Davidson County, but is separated by one tract in between K3 observations. Second-order queens contiguity resulted in a p-values <.05 across all census tracts. This finding leads to a rejection of the null hypothesis that the observed pattern was due to chance and may suggest a non-random neighborhood shaping phenomena.

Table 4. *Neighborhood change typologies and associated join-count statistic p-values.*

	First-Order QC p-value	Second-Order QC p-value
K1	0.0001	0.0001
K2	0.0001	0.0001
K3	0.0226	0.0417
K4	0.5995	0.0065

Supervised Machine Learning

Supervised machine learning uses input variables (Y) to estimate a mapping function that converges on the value or class of the output variable (X). Labeled outcomes associated with the input variables, allow for model performance validation. In short, predictive supervised models learn from training data, are assessed using testing data, and ideally perform well enough to be extrapolated to future data. The concept of generalization describes the degree to which a machine learning model has correctly interpreted information from the training data to cases the model has not encountered during the training phase. A useful model is able generalize from inductive learning during the training phase to make correct predictions when introduced to an unseen dataset.

We use the outcomes and interpretations of the clustering algorithm to label census tracts as gentrifying or not gentrifying. We identified one neighborhood change typology (k2) as exhibiting the characteristics associated with gentrification. For the purposes of our predictive model, the remaining neighborhood typologies are collapsed into a single category of “not gentrifying”. We then employ two separate supervised machine learning algorithms to construct predictive models capable of identifying gentrification from base year predictive variables (table x). Principle component logistic regression (appendix) and random forest algorithms were chosen based on predictive power and interpretability.

Random Forest

Random forest is an ensemble machine learning method which begins by bootstrap aggregating, or bagging, to partition a dataset into many random subsets of training and testing data. Bagging is fundamental to the random forest algorithm and has the effect of reducing prediction variance. Bagging subdivides the entire population of observations into many different training and testing sets, builds predictions on these iterations, and then averages the predictions to achieve a lower variance. Bagging is used to grow many deep and unpruned trees. A single tree suffers from high variance, but low bias. Averaged across thousands of iterations however, the prediction converges to a lower variance and typically improves accuracy. As an ensemble learner, the random forest algorithm applies the “wisdom of the crowd” to classification problems by taking a majority vote of the most commonly occurring class from the crowd-sourced population of decision trees.

Bagged training data is used to construct individual classification trees while the testing, or out of bag samples, are used to assess model performance. A random subset of predictive variables are then considered as possible branches at decision nodes. The feature-value combination of only the randomly selected subset of predictive variables that maximizes information gain or decreases gini impurity⁷ is selected, resulting in a binary split with more homogenous, or “pure”, downstream partitioned observations (Breiman, 2001). This randomization of available variables in each node helps to reduce the correlation between trees and improves predictive power and classification accuracy. Each decision tree partitions observations until a classification is achieved at the point when no further class purity can be achieved. Individual decision-tree classifications are then voted on towards the final predictive classification with a lower variance and higher predictive power than any single tree alone (Mellor et al., 2013). The out-of-bag approach was used to assess model performance from observations not involved in the bootstrapping model tuning process. Random forest classification was chosen above other model types because of its generalizability, ease of tuning, and insensitivity to outliers.

Table 5. Random Forest Predictive Variables (2000).

Demographics	Housing	Occupational	Transportation	Amenities
% Poverty	% Vacant Land	% Unemployed	% Public Transit Commute	% Park coverage
% Non-Family	% Dwelling 5+ Units	% Blue Collar	% No Car Available	Distance to Downtown
% Aged Under 18	% Renters		% 3+ Cars Available	
% Aged Over 65	% Rent Burdened		% Commutes less than 20 min.	
	Eviction Rate			
	Housing Stock Age			

⁷ Gini impurity is a measurement assessing the likelihood of a randomly chosen nodal observation being mislabeled from the distribution of responses in the node.

Prediction Performance

We tested the predictive power of random forest classification using a training set (2000 – 2016) of 70%, or 107 observations, of observations and test the model performance on the remaining 30% (45 observations). We employ the R package ‘caret’ to streamline data splitting, model tuning, and validation associated with our predictive models (Kuhn, 2019). The overall class balance of 13% positive (gentrifying) tracts was maintained proportionally between training and testing data sets to ensure that each set was representative of the population. The models are trained and tuned using the training set and are then evaluated against the unseen test data. Leave one out cross validation was used to tune the model on the training datasets and to maximize the limited sample size.

We evaluated the performance of the random forest (RF) algorithm’s ability to classify gentrifying tracts apart from non-gentrifying tracts using data from the base year. We selected a classification criteria that balanced specificity and sensitivity (.49) rather than accuracy due to the class imbalance nature of the problem. The tradeoff between sensitivity and specificity via changing classification thresholds can be observed in the ROC plot. Figure 18 shows the RF model trained perfectly on the training set. While initially this may appear indicative of overfitting, tuning measures, RF statistical properties, and the validation/test set approach all explain or guard against model overfitting.

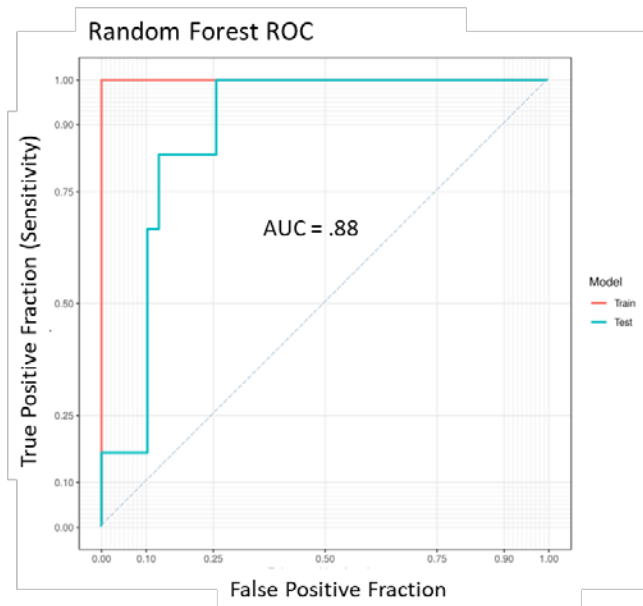


Figure 18. Random Forest trained model Receiver Operating Characteristic (ROC) curve.

Table 7 provides model summary statistics for our RF model, many of which can be derived from the confusion matrix in table 6. Because of the class imbalance and low percentage of true positive

classes, some of the metrics are more representative of model performance than traditional metrics such as accuracy. Additionally, a consideration of the problem context will help evaluate the usefulness of the model. Similar to fraud or medical diagnoses, there may be few positive cases, but a high cost associated with a false negative prediction. In this context, that is predicting a census tract will not gentrify when it actually does gentrify. In such cases where the cost of false negatives is high, recall/sensitivity is a valuable metric. Sensitivity is defined as the predicted positive cases divided by the total positive cases. At a threshold value of .49, our RF model correctly identifies five out of six (83%) of gentrifying tracts. This is associated with a precision of 53% of our predicted positive cases. The F1 score is another intra-model performance metric derived from the harmonic average of the precision and recall, which yields .55 in our selected model.

Table 6. Random forest model confusion matrix.

Predicted	Observed	
	Not Gentrifying	Gentrifying
Not Gentrifying	33	1
Gentrifying	6	5

Table 7. Random forest model performance statistics.

Metric	Value
Accuracy	0.844
Kappa	0.502
Sensitivity	0.833
Specificity	0.846
Precision	0.455
F1	0.588
Balanced Accuracy	0.840
AUC	0.885

Other metrics like balanced accuracy and kappa are useful to evaluate overall model generalizability (Brodersen et al. 2010;). Balanced accuracy attempts to minimize the bias towards the more frequent class by averaging both class accuracies, giving $.5 * \left(\frac{TP}{P} + \frac{TN}{N} \right)$. The RF model balanced accuracy resulted in a value of .84. This value is close to the standard accuracy metric of .83 suggesting that the model is performing similarly for both positive and negative classes. Additionally, the Kappa statistic considers observed accuracy with respect to expected accuracy due to random chance. Our RF Kappa value of .5 falls into the “moderate agreement” that the observed performance is performing better than chance (Landis and Koch, 1977; Viera and Garrett, 2005).

Variable importance is reported by the mean decrease in the Gini index in table 8. This metric reports each predictive variable’s splitting power of homogenous child nodes. An important feature is

one that splits the data into homogenous or pure classes (Breiman, 2001). The RF model learned heavily from labeled data that the majority of gentrifying tracts were close to downtown. Additionally, socio-economic variables such as poverty, and vacancy, and unemployment were shown to be informative variables regarding the predicted outcome of gentrification. Less important variables were blue collar workers, evictions, and renting dynamics. Built environment variables such as new housing and housing with 5+ units fell in the middle of the pack, while park coverage was the sixth most important variable considered. Because predictive variables are correlated, variable importance measures may be unstable and vary from run to run. This effect is also a relic of a small sample size of the data set. Nevertheless, variable importance can provide valuable insight into the starting characteristics of an area that may be advantageous to gentrification.

Table 8. *Random forest variable importance.*

Ranking	Variable	Importance
1	DOWNTOWNDISTANCE	100.00
2	PCTPOV	48.78
3	PCTVACANT	35.74
4	PCTUNEMPLOY	31.69
5	UNDER18	29.66
6	PARKCOVERAGE	23.23
7	PCTBIGUNIT	21.47
8	NEWHOUSING	21.36
9	PCT3CAR	20.55
10	PCTNONFAM	18.32
11	PCT20MINCOMMUTE	18.03
12	OVER65	15.78
13	PCTRENTENTER	14.34
14	PCTRENTBURDEN	6.12
15	EVICTONRATE	1.61
16	PCTBLUECOLLAR	0.00

The small sample size of the observations used to build the predictive model resulted in prediction variation between model runs. Some variation was expected because of the inherent stochasticity in the random forest algorithm as well as the random selection of observations used to build the testing and training data sets. However, a sensitivity analysis using a Monte Carlo simulation quantified the variation in the most probable tracts predicted to gentrify. Figure 19 shows the variation between model runs for each census tract. The median prediction value of 200 separate model runs is plotted on the x-axis against the range of observed values on the y axis.

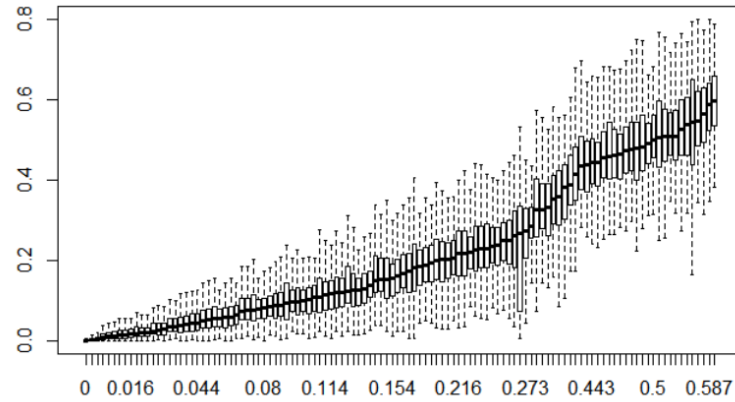


Figure 19. Random forest model prediction variability. Median prediction values (x-axis) plotted against 200 model run values.

Mapped Predictions

The previous section outlined the results on holdout set to evaluate the generalizability of the predictive algorithm. Here we apply the model to unseen 2016 data to predict future gentrification susceptibilities. Figure 20 maps the corresponding potential for future gentrification based on the RF model trained on 2000 data. The symbology and associated probabilities are grouped broadly due to variations between model runs that will be discussed further in the following section.

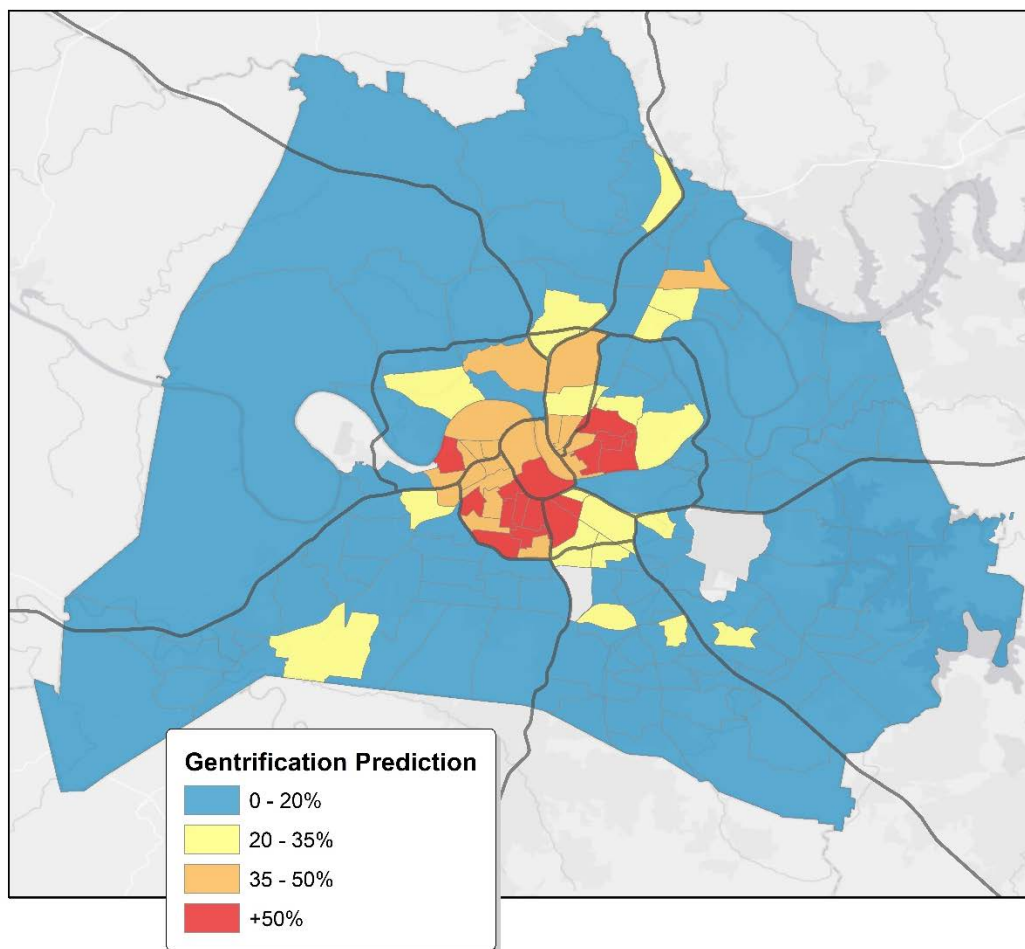


Figure 20. Random forest model future predictions using 2016 data.

There are three arms of high-risk areas in south, east, and northwest Nashville. A larger scale map of the inner-city area is provided in figure 21. The first high probability pocket is apparent in south Nashville, situated in between the I-40 and I-440 highways. Three of the four tracts in this area that are greater than 50% were previously identified as gentrifying between the 2000-2016 time period. These neighborhoods include Edgehill, Wedgewood-Houston, the Fort Negley area, and Chestnut Hill.

The northern arm of high risk areas lies immediately north of I-40/I-65 and bounded by the Cumberland River. One census tract north of I-40 and I-65 is predicted as the highest (>50%) risk areas for future gentrification. This census tract is neighbored to the east by several moderately predicted tracts as well. This area represent a shift further north than previously identified gentrified areas, which were identified south of the interstates. Neighborhoods in these areas include Buena Vista Heights, Cumberland Gardens, College Heights/Clifton, and the area surrounding Tennessee State University.

The eastern arm of high risk areas is located immediately east of the Cumberland River. Four census tracts in East Nashville fall into the highest risk areas- all of which were identified as previously gentrifying and are predicted to continue to gentrify. These neighborhoods include Shelby Hills, Five Points, and areas surrounding McFerrin and Cleveland Park.

Other noteworthy spatial trends from the RF predictions show Interstate 440 as a boundary between many moderate risk (35-50%) tracts apart from the lowest vulnerable tracts to the southwest. The Cumberland river also differentiates many of the moderate risk areas in North Nashville from less at-risk areas north of the river. The mapped RF predictions also reveal interesting patterns further away from the urban center. There appears two corridors to the northeast and to the southeast, both bounded by I-24 and I-65 with risk predictions in the range of 20% - 35%. Although only modest risk predictions, these areas were identified as previously affordable areas where those displaced from gentrification may have resettled.

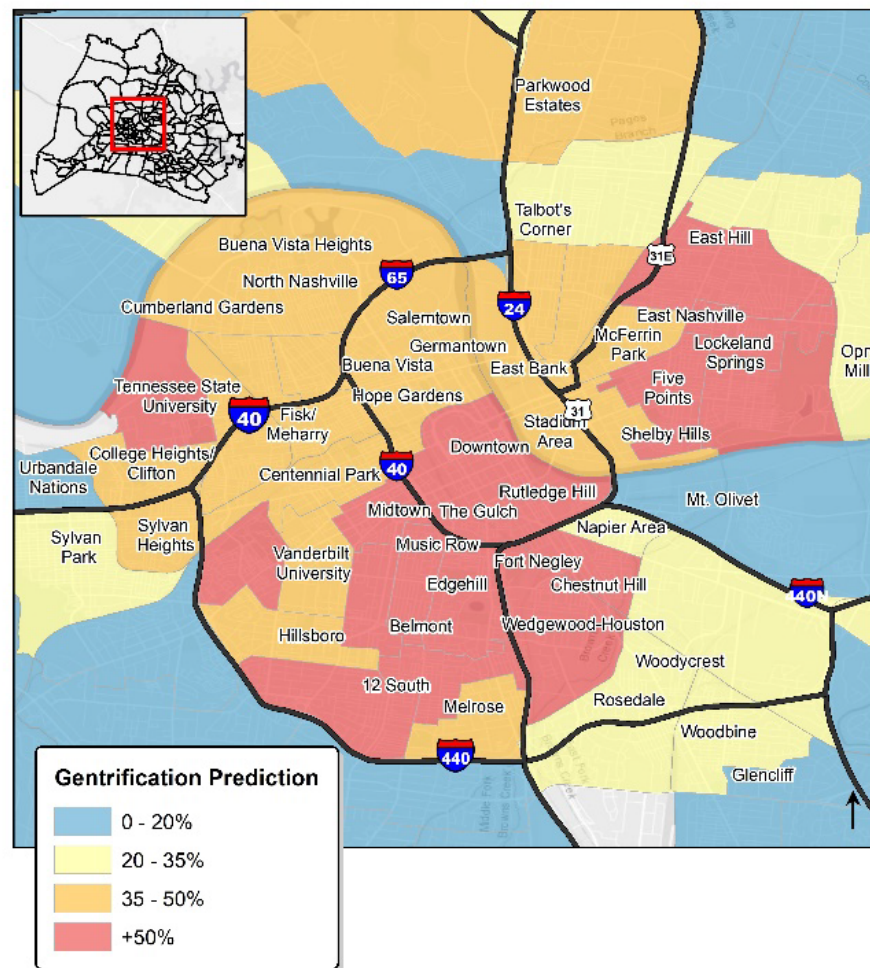


Figure 21. Zoomed gentrification predictions.

Teardowns

We evaluated the agreement of our gentrification predictions against the location of residential teardowns in Davidson County. Teardowns are an improvement to the built environment where the existing residential structure is completely demolished with the immediate intention of replacement by a new housing unit (Munneke and Womack, 2013; Weber et al., 2006). The decision to tear down a habitable housing unit is directly explained by Neil Smith's rent gap theory. As the land value of a home grows irrespective of the deterioration of the physical property, demolition and redevelopment represent a rational market response to capitalize on the untapped ground rent (Smith, 1996).

Teardowns are one of the most visible indicators of the built environment effects of gentrification. They simultaneously exploit and contribute to rising neighborhood land values (Munneke and Womack 2013). Teardowns have become a familiar sight in urban gentrifying neighborhoods; they represent a multi-faceted proxy for gentrification and displacement as well as the social and financial pressures on incumbent residents.

However not all gentrifying neighborhoods may experience residential teardowns and not all teardowns occur in gentrifying neighborhoods. Weber et al. (2006) examined the different physical, economic, political, and cultural variables that explain demolition activities in Chicago. They found that built environment characteristics were the strongest predictors of demolition. Specifically, smaller, older frame homes with less lot coverage were more vulnerable to demolition. Political jurisdictions and community demographics like race explained little of the variation between demolition susceptibility. Locational characteristics like distance to tax increment financing (TIF) districts suggest that government sponsored intervention strategies may have an impact in suppressing teardown redevelopment (Weber et al., 2006).

Residential teardowns were identified using an open-source building permit database (Metro Codes Department, 2018). Teardowns were operationally defined as addresses with both a demolition permit followed by a new residential construction permit. Accessory structures like garages, carports, and pools were excluded from the analysis. We identified 771 residential teardowns between 1/1/2016 and 4/2/2018. Figures 22 and 23 illustrate two examples of residential teardowns. A street-level inspection of addresses identified as teardowns confirmed that these homes were either newly constructed, under construction, or demolished and pending construction.

June 2016



November 2018

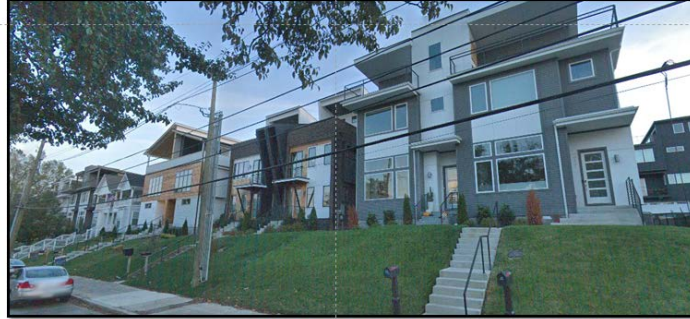


Figure 22. South Street in the Nations neighborhood (Google Street View).

June 2016



March 2017



Figure 23. Residential teardown in the East Nashville neighborhood (Google Street View).

Figure 24 shows the location and construction cost associated with the residential teardowns as well as our gentrification risk predictions. Teardowns exhibit both spatially clustering at the neighborhood level as well homogeneity in the construction cost of the new structure. There appears moderate agreement between our risk predictions and the location of residential teardowns with two distinct exceptions.

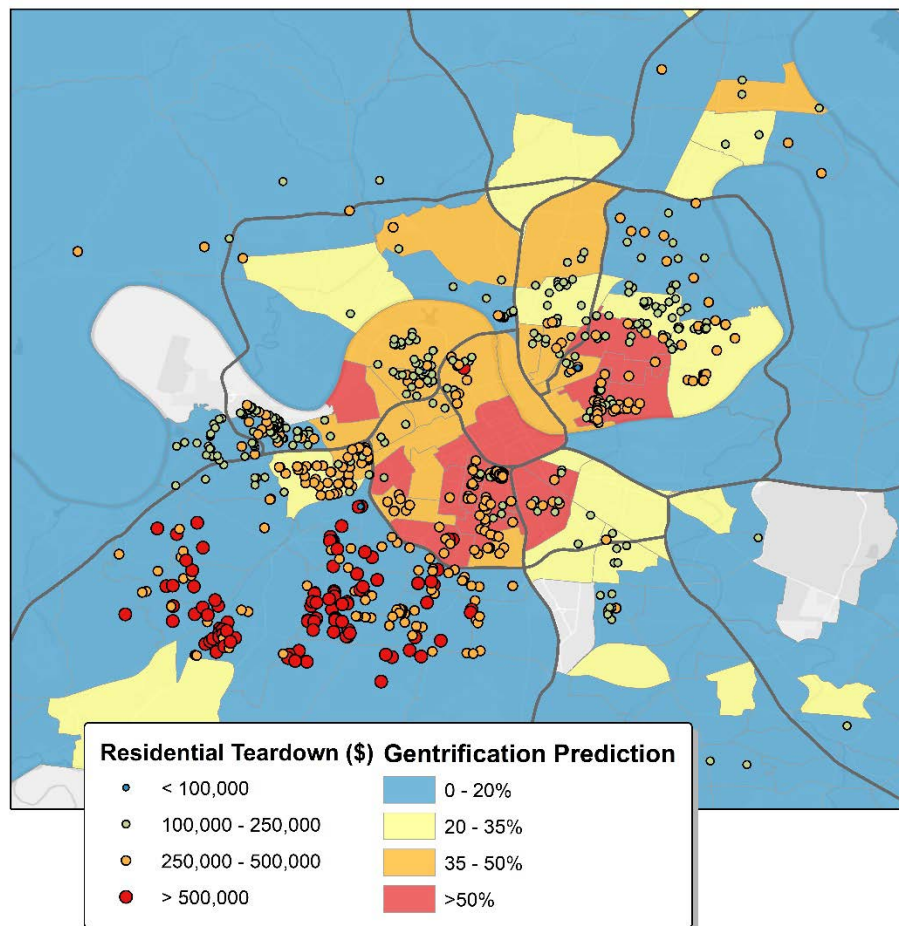


Figure 24. Gentrification predictions (2016 data) plotted against residential teardown locations and new construction values (2016 - 2019) (Metro Codes Department, 2019).

Teardowns with the highest construction cost (greater than \$500,000) are highly spatially clustered in the southeast portion of the county where home prices are the highest. Podagrosi et al. (2011) found a similar upgrading process in a wealthy Houston neighborhood, Bunker Hill Village. They identified large, expensive, non-deteriorated homes that were increasingly being demolished and replaced by even larger, more expensive, customized homes. Nashville's wealthiest neighborhoods like Belle Meade and Green Hills may be experiencing similar upgrading processes as those found in Houston.

The other noteworthy exception between the location of residential teardowns and our predictions is located in the Urbandale Nations neighborhood. There are over 70 teardowns clustered in this neighborhood, which is located north of I-40 and south of the Cumberland River. This area has been anecdotally linked to gentrification, yet in our analysis, it was not identified as gentrifying between 2000 and 2016 (Ward and Reicher, 2017) and is not expected to gentrify. The Nations is certainly experiencing an upgrading process of some kind, however it may be more in line with revitalization than gentrification. Between 2000 and 2016 home values increased by 73%, but rent only increased by 8%. Additionally, incomes grew by only 4%, the percentage of non-white residents grew by over 6%, and multi-unit housing increased by nearly 8%. The small rent and income increases, expansion of multi-unit housing, and increasing racial diversity were enough to differentiate the neighborhood apart from the gentrification cluster which experienced more drastic rent and income increases, as well as a loss in non-white residents and multi-unit structures. These discrepancies led to the area's low gentrification probability, but the appearance of a significant cluster of residential teardowns as it experiences improvements of the built environment without a substantial change to the community demographics.

Discussion

Identifying Gentrification

This first component of this research advances an alternative methodology to identify gentrification in Davidson County between 2000 and 2016. We identified four distinct neighborhood change typologies via a k-means clustering procedure on percent changes of six socio-economically important census variables. One specific typology exhibited the distinct markers of gentrification, signaled by the largest increases in home value, income, and educational attainment. This typology also revealed the largest decrease in non-white persons across any typology. It should be noted that these are relative measurements rather than net mobility, but is highly suggestive of residential displacement within gentrifying areas. This evidence for displacement was used to differentiate the process apart from revitalization or incumbent upgrading.

In total, 20 census tracts, or 13 percent of all Davidson County tracts, were distinguished as gentrified between 2000 and 2016. Compared with previous methods to identify gentrification in Davidson County, the clustering results here identified fewer tracts than the Freeman method (37), and more than either the Governing (12), Ellen O'Regan (8) or McKinish et al. (5) methods⁸. The spatial clustering of gentrifying areas around the urban center meets the expectation that the process may be limited to center-city tracts without explicit spatial constraints. These locations corroborate anecdotal accounts of gentrification within neighborhoods such as Germantown, Edgehill, and East Nashville (Plazas, 2017).

⁸ Measured from 2000 – 2010.

K-means clustering provided a more holistic representation of neighborhood change than prescriptive efforts to distinguish the process. In addition to identifying gentrifying areas, the method revealed areas of affordability in the northeastern and southeastern corridors of Davidson County. These areas are a striking inverse of the gentrification profile and potential locations that those displaced from gentrifying neighborhoods could relocate. The affordable typology showed the largest decreases in home value, rent value, and income. As such, these areas are likely experiencing an increased concentration of poverty. This typology also contained the largest increase in multi-unit housing as well as the largest increase in non-white persons. These areas are found further from the urban core of Nashville, but concentrated along the major northeastern and southeastern interstate arteries. The spatial patterning of these potential relocation areas to away from the inner-city is consistent with displacement relocations found in New York City (Wyly et al., 2010). This finding is also in line with previous research that suggests accessibility to labor markets via highways may be a factor in relocation housing preferences (Voith, 1993).

The spatial clustering of all four neighborhood change typologies proved to be significant- tracts have a higher probability of exhibiting the same neighborhood change typology as their adjacent tracts than random chance. Intuitively, this finding suggests that signals of investment and disinvestment in one area may impact perceptions on adjacent neighborhoods in a reinforcing role. The spatial dependence of neighborhood trajectories has been observed in other cities including Chicago, Detroit, and Phoenix (Ling and Delmelle, 2016). Brown-Saracino (2017) describe gentrification as “contagious”, suggesting its uneven presence in areas that are located nearby other gentrifying areas. The spatial structure of neighborhood change provides further evidence that gentrification is a selective process rather than a sporadic one. In turn, this provides motivation for the potential to predict gentrification before it takes hold in a neighborhood. Furthermore, the spatial information uncovered here demonstrates the importance of endogenously engineered variables that factor in adjacent characteristics. This extends the scope of the gentrification process beyond a single census tract, but importantly considers neighboring factors into future risk considerations.

The machine-clustered gentrification typology is also tailored to Nashville, Tennessee. Shaw (2005) notes that gentrification “plays out differently in different places and the process is deeply affected by the local context.” As such, we would expect to find different paces, magnitudes, and typologies of change between Nashville, New York, San Francisco, or Atlanta. The method outlined here affords these contextual differences, while also allowing for the cross-comparison of cities “gentrification fingerprints.” The Nashville gentrification fingerprint is marked by large increases in home/rent values, incomes, and education with simultaneous large decreases in non-white persons. The clustering approach used in this study provides this transparent, city-specific quantitatively reportable definition for gentrification without the dictatorial frameworks of previous research. Future work would benefit from evaluating both the spatial extent and magnitudes of the gentrification profile across different cities.

This method also has the capacity to quantitatively distinguish neighborhood revitalization apart from gentrification. Local policy-makers may uncover the factors that support these more equitable outcomes in order to reverse engineer successful intervention strategies and optimize resource allocations.

Lastly, the clustering approach offered here may be refined with different input variables that capture investment, disinvestment, and demographic markers of gentrification. City-specific data

sources like building permits, evictions, code violations, and home appraisal/sales values are just some of the untapped metrics that can be incorporated into this multi-dimensional framework.

Predicting Gentrification

Over the course of the past 50 years perceptions of gentrification have evolved from a sporadic, localized process to a pervasive threat to disadvantaged communities across the globe. During this same period, however, gentrification research has been primarily backward looking and concerned with causes and effects. This research has, however, advanced the empowering notion that gentrification is a rational response to changing market conditions and consumer preferences. This idea is inherently valuable because it implies that gentrification may be predictable to some reliable, or useful degree. Neighborhoods downgrade and upgrade according to differentiable trajectories that are a function of both internal and external factors. We set forth a machine learning framework that attempts to capture these factors relative to historical patterns of neighborhood change. We ask if there are discernable housing, demographic, occupation, transportation, amenity, or accessibility characteristics that may explain why certain neighborhoods are more likely to gentrify.

The second component of this research demonstrated the performance of supervised machine learning framework to identify at-risk areas for future gentrification. This method used the gentrification typology from the k-means clustering process to label past observations and then deduce the starting characteristics that may influence the observed outcomes. Therefore, the predictions represent the probability of gentrification based on its observed character in the same area, which are marked by large increases to rent, home, income and education, as well as large decreases in non-white residents and multi-unit housing. We feel confident that the local magnitudes of these changes that were observed in the past provide the best indicator of how gentrification will manifest in the future.

We collapsed the gentrification typology against all other neighborhood change typologies. From this, we trained a random forest classification model to learn from a variety of housing, family structure, transportation, and built environment conditions. A validation set approach was used to evaluate the performance of the model and determine its generalizability on future, unseen data. The random forest model correctly identified 83% of gentrifying tracts at the expense of a 45% false-positive rate. In the context of gentrification, where the possibility of overlooking a gentrifying tract is costly (type II error), the high sensitivity of the model provided confidence that the model had identified useful and distinguishable patterns that could be projected on 2016 data.

The most important feature used by the random forest model was distance to downtown followed by percent poverty, vacancy, and unemployment. Distance to downtown was by far the most significant variable in our model, however other researchers like Chapple (2009) found distance to San Francisco was the 16th (out of 19) most significant predictor of gentrification. Distance was calculated geodesically, or “as the crow flies”, from the census tract centroid to downtown. However other distance measures may provide additional valuable predictive insight like road network distance and/or travel times during peak demand hours.

Several of the less important features also ranked relatively low on previous studies attempting to identify the conditions behind gentrification. Chapple (2009) similarly found rent burdens (11th), non-family households (8th) ranked as relatively low factors. Percent renters ranked surprisingly low in feature importance. We expected the renter-homeowner dynamic to be more prevalent with areas featuring larger proportions of renters as potential paths of least resistance for evictions and other mechanisms preceding gentrification. Additionally, eviction patterns did not appear to signal future gentrification. Displacement pressures within Nashville may be attributed more to social isolation and exclusion rather than more formal economic displacement mechanisms.

The spatial distribution of the gentrification risk predictions revealed two distinct arms of the highest classified risk (>50%) category. These high risk areas reaching to the south and east contained nine census tracts that were previously identified as gentrifying between 2000 and 2016. The predictions suggest that these areas may be expected to continue to rapidly gentrify based on their proximity to downtown, remaining economically vulnerable populations, and vacancy characteristics. Lees (2003) coined the term “super-gentrification” to describe an intensified form of mature gentrification that can transform a working class neighborhood into an elite enclave.

A third area with moderately high risk (35 – 50%) lies to the north of I-440/65 loop and constrained to the south by the Cumberland River. The majority of areas in the North Nashville area were not previously identified as gentrifying. These areas may represent attractive locations for redevelopment because of their close proximity to academic institutions that include Tennessee State University, Fisk University, and Meharry Medical College. The presence of higher education campuses was not explicitly coded for within the model, but has been theorized to attract higher-paid “creative professionals” who value diversity and community above traditional amenities like shopping malls and sports stadiums (Florida, 2003; Rutheiser, 2011; Weissbourd et al., 2009). Additionally, we identified areas adjacent to Vanderbilt University and Belmont College that experienced gentrification between 2000 and 2016. The cultural diversity, talent, intellectual resources of these institutions may be attractive housing anchors for the highly educated portion of the workforce that is often associated with gentrification. This may be even more pronounced from broader economic shifts away from an industrial-based economy to a knowledge-based economy that increasingly values human and intellectual capital (Florida and Cohen, 1999).

Another finding of the random forest predictions suggests that the process may be expected to spill over into previously affordable areas further away from the downtown core. Despite the distance to downtown variable’s high feature importance within the model, there appears two corridors of moderate (20-35%) risk extending towards the northeast and southeast. These corridors are bounded by I-24 and I-65 and coincide with the affordability typology- areas where rents and home values lagged behind the Davidson County medians. We interpreted these areas as beacons of affordability undergoing significant demographic change, likely as a result of the relocation of those displaced within gentrifying areas. Although the predictions are on the lower end, this is a potentially troubling signal suggesting the threat of gentrification and displacement could encroach again on those that were previously displaced to these areas. This process is supported by Neil Smith’s geographic explanation for investment decisions where focused investment in one area leads to market barriers in that area as well as underdevelopment in other areas (Smith, 1979). The previously disinvested areas begin to look more attractive to developers as the rent gap widens and the potential for profit increases. Despite the distance to downtown variable’s high feature importance within the model, there may be some of the

same demographic vulnerability indicators like high poverty and vacant land in these areas that spurred gentrification in the past. This finding is also consistent with Butler (2007), who similarly acknowledged gentrification spreading beyond the central city limits.

The predictive component of this research advances the potential for more forward-looking research in the context of gentrification. Mitigative policy interventions designed to protect and expand access to affordable housing may find more success than adaptive strategies that respond only during or after gentrification has a foothold on a neighborhood. Over 50 years of research on gentrification has advanced our understanding of the theoretical causes of gentrification, as well as the specific features like accessibility, public transit, amenities, and creative centers that drive consumer preference and profitability. Economic theories advanced by early researchers suggest that gentrification is as a rational consequence of supply and demand side in disinvested areas. These advancements bolster the primary assumption of this research component- gentrification can be predicted. In addition to theoretical advancements we have decades of evidence of certain areas that have gentrified and areas that have not. This historical evidence can be put to work in a supervised machine learning framework that learns from a city's history.

Limitations and Future Directions

We treat neighborhood change as discrete outcomes between periods in time. An argument may be made to analyze shorter intervals over longer time periods. Previous empirical research on gentrification has largely been constrained to decadal periods of analysis due to the frequency of the US Decennial Census. The American Community Survey (ACS) replaced the Decennial Long-form census in 2010. The ACS captures shorter time intervals allowing for more temporally fine-grained resolution of neighborhood change at the expense of a smaller sample size. The ACS allows for shorter interval change “snapshots” between two points in time, albeit at a smaller sampling size with accuracy concerns (Bazuin and Fraser, 2013). A similar analysis could use shorter intervals to build a larger training and testing dataset than that used in this study.

Our predictive model does not take into account certain important, but unforeseeable market disruptors like housing market crashes, the adoption of new public transportation systems, natural disasters, or major policy changes. The model instead relies heavily on the rationality of profit-seeking developers as well as certain amenity features that may be disproportionately attractive to individuals who are able to pay a premium to access these services in previously disinvested areas.

We initially set out to use a variety of open source data specific to Nashville as predictive variables for our predictive model. These features included residential teardowns, renovations, building quality, code violations, business locations, and police service calls. Although these data sources are openly available, they are often limited to only a few years in scope. We were ultimately limited to data that was available in both 2000 and in 2016. Business locations could potentially be a valuable predictive variable, as the influx of higher income serving amenities like coffee shops and boutiques signal a form of retail upgrading. However, data obtained from the Yelp Fusion API did not include information as to when the business was established. Resolving this matter may require the cross-

comparison of business permits, local administrative data, and/or non-open source data, but could provide valuable insight into commercial and retail influences on gentrification. We encourage the expansion and refinement of both the demographic and investment variables used in the k-means clustering identification method as well as the predictive variables.

The temporal limitations of using 2000 data to build our predictive model also influenced the use of our spatial units of analysis. Because the US Census tract boundaries have changed over time, a normalized dataset to 2010 boundaries was required. We used the Geolytics neighborhood change database- a commonly used dataset in neighborhood change analyses. Smaller spatial units like census blocks and block groups are not available through this resource. These smaller units would provide a higher spatial resolution of neighborhood change; the average size of a census tract in Davidson County is approximately 3.3 square miles compared to 1.1 square miles of the average block group. A census block-based analysis would also increase the sample size from 161 total census tracts to 473 block groups, elucidating several of issues of sample size in the predictions.

The small sample size of the observations used to build the predictive model resulted in reported variation between model runs. Variation is expected to some degree because of the inherent stochasticity in the random forest algorithm as well as the random selection of observations used to build the testing and training data sets. However, a sensitivity analysis using a Monte Carlo simulation quantified the variation in the most probable tracts predicted to gentrify. Increasing the sample size would likely reduce the observed variability. To offset this effect, we presented the median predictions out of 200 model runs.

Lastly, researchers relinquish that the gentrification process can take many different forms depending on the city and time period, yet they continue to conduct nation-level generalizations. We advocate for more city-specific and time sensitive research on gentrification. Kennedy and Leonard (2001) state “The quickly changing nature of forces driving gentrification conflicts with the methodical pace of bureaucracies and the long timeframes required by many of the financing-and construction-based strategies needed to address it.” In a similar vein, we should not rely on outdated research to guide our current views on the process as we have seen it evolve in considerable magnitude and effects since its inception. Ultimately, academic efforts designed to improve the social equity in cities. We propose that this method be considered on a city by city basis in order to capture the city-specific character of emergent trajectories of neighborhood change.

Conclusion

Throughout this paper we have encouraged a back to the basics approach to defining and assessing gentrification as well as a critical examination of previous operational definitions of gentrification. Identifying gentrification is a high-stakes classification exercise; Classification is in its purest sense a form of simplification. Any effort to classify and label gentrification stem from our human propensity to categorize the world around us. This penchant has been beneficial to human

development, helping us to make sense of new observations. However, this ingrained habit to simplify, classify, and categorize carries important consequences within the current context. The radical diversity of methodologies used by previous researchers to classify gentrification is a closer reflection of individual experience, perceptions, personal heuristics, and expectations than it is of objective science. To this point, urban change research is largely responsible for crystallizing public and policy-maker perceptions on gentrification. The wealth of contradictory evidence, largely driven by inconsistent methodologies, has permitted both supporters and opponents of gentrification a means to selectively entrench their views on the subject. Gentrification research may benefit from the emotionless 21st century tools available that rightfully resituate the emphasis back towards the underlying data patterns and away from the dictatorial and prescriptive frameworks of the past.

We also must explicitly consider the diversity of actors involved in gentrification and recognize the financial incentives and profit that are at stake. Wyly et al. (2010) note that “one of the most effective tactics of neoliberalism involves the statistical disappearance of its costs and victims.” An inconsistent lexicon surrounding gentrification has fostered a prescriptive class of measurements that place an undue amount of influence through arbitrary thresholds. These manufactured frameworks easily stray from their intended target, and serve as control knobs with the potential to severely bias results. The pairing of a tool as unemotional as machine learning with a deeply personal experience like gentrification feels contradictory. However, it is precisely because of this dichotomy that we must explore more objective tools that identify real patterns rather than reproduce assumptions and biases. As the conceptual understanding of gentrification has evolved, the techniques used to measure the process should similarly evolve. Researchers have the responsibility to explore less biased alternatives that are capable of capturing the multiple dimensions and complexity that the process deserves. Consistent, transparent, and contemporary will go a long way in our efforts to create more socially just and equitable communities.

Lastly, gentrification research may find more success with forward-looking research that are mitigative rather than adaptive in focus. Returning to Ruth Glass’s (1964) original definition of gentrification, “Once this process of 'gentrification' starts in a district it goes on rapidly until all or most of the working class occupiers are displaced and the whole social character of the district is changed.” Policies to combat gentrification may be too late to respond to gentrification once it is observed on the ground level. Reactionary policies have proven largely inadequate in their attempts to stall the economic and preference-based engines powering gentrification. However, policies that are implemented before gentrification takes a hold of a neighborhood may find more support. The machine learning framework applied here puts decades of past observations to work in order to identify the underlying characteristics that may lead to selective patterns of gentrification. Cities may better position their limited resources in these areas to promote more equitable development. These policies take the form of rental housing demonstration (RAD), housing choice vouchers, tax credits, or inclusionary zoning. Regardless of the policy instrument used, they should broadly focus on building existing residents’ capacity to stay. Contrary to some beliefs, this is not mutually exclusive to growth and development, but certainly involves the preservation and expansion of affordable housing options for those lacking a political voice.

References

- Airbnb, Housing, and Nashville, 2017. . Airbnb.
- Allison, M., 2017. Nashville Tops the List of Hottest Housing Markets for 2017. Zillow Porc. URL <https://www.zillow.com/blog/hottest-housing-markets-2017-209986/> (accessed 6.16.19).
- Anselin, L., 2018. Contiguity-Based Spatial Weights [WWW Document]. GeoDa. URL https://geodacenter.github.io/workbook/4a_contig_weights/lab4a.html#queen-contiguity (accessed 6.23.19).
- Armstrong Jr, R., 1994. Impacts of Commuter Rail Service As Reflected in Single-Family Residential Property Values. *Transp. Res. Rec.* 88–98.
- Arthur, O., 2005. Gentrification and crime. *J. Urban Econ.* 57, 73–85.
- Atkinson, R., 2004. The evidence on the impact of gentrification: new lessons for the urban renaissance? *Eur. J. Hous. Policy* 4, 107–131.
- Atkinson, R., 2000. Measuring Gentrification and Displacement in Greater London. *Urban Stud.* 37, 149–165. <https://doi.org/10.1080/0042098002339>
- Barton, M., 2014. An exploration of the importance of the strategy used to identify gentrification. *Urban Stud.* 53, 92–111. <https://doi.org/10.1177/0042098014561723>
- Bazuin, J.T., Fraser, J.C., 2013. How the ACS gets it wrong: The story of the American Community Survey and a small, inner city neighborhood. *Appl. Geogr.* 45, 292–302. <https://doi.org/10.1016/j.apgeog.2013.08.013>
- Berry, B.J.L., 1982. *Islands of Renewal, Seas of Decay: The Evidence on Inner-city Gentrification*. School of Urban and Public Affairs, Carnegie-Mellon University.
- Bivand, R., 2019. *spdep: Spatial Dependence: Weighting Schemes, Statistics and Models*.
- Bostic, R.W., Martin, R.W., 2003. Black Home-owners as a Gentrifying Force? Neighbourhood Dynamics in the Context of Minority Home-ownership. *Urban Stud.* 40, 2427–2449. <https://doi.org/10.1080/0042098032000136147>
- Bowes, D.R., Ihlanfeldt, K.R., 2001. Identifying the Impacts of Rail Transit Stations on Residential Property Values. *J. Urban Econ.* 50, 1–25. <https://doi.org/10.1006/juec.2001.2214>
- Breiman, L., 2001. Random Forests. *Mach. Learn.* 45, 5–32. <https://doi.org/10.1023/A:1010933404324>
- Brown-Saracino, J., 2017. Explicating Divided Approaches to Gentrification and Growing Income Inequality. *Annu. Rev. Sociol.* 43, 515–539. <https://doi.org/10.1146/annurev-soc-060116-053427>
- Butler, T., 2007. For Gentrification? *Environ. Plan. Econ. Space* 39, 162–181. <https://doi.org/10.1068/a38472>
- Cervero, R., Duncan, M., 2004. Neighbourhood Composition and Residential Land Prices: Does Exclusion Raise or Lower Values? *Urban Stud.* 41, 299–315. <https://doi.org/10.1080/0042098032000165262>

- Cervero, R., Landis, J., 1997. Twenty years of the Bay Area Rapid Transit system: Land use and development impacts. *Transp. Res. Part Policy Pract.* 31, 309–333.
[https://doi.org/10.1016/S0965-8564\(96\)00027-4](https://doi.org/10.1016/S0965-8564(96)00027-4)
- Chapple, K., 2009. Mapping Susceptibility to Gentrification: The Early Warning Toolkit. University of California.
- Chatman, D.G., Tulach, N.K., Kim, K., 2012. Evaluating the Economic Impacts of Light Rail by Measuring Home Appreciation: A First Look at New Jersey's River Line. *Urban Stud.* 49, 467–487.
<https://doi.org/10.1177/0042098011404933>
- Clarivate Analytics, 2019. Web of Science [WWW Document]. URL https://wcs-webofknowledge-com.proxy.library.vanderbilt.edu/RA/analyze.do?product=WOS&SID=5BfsL5Q15wJ8dwbN9wz&field=PY_PublicationYear_PublicationYear_en&yearSort=true (accessed 6.18.19).
- Clark, E., 1995. The Rent Gap Re-examined. *Urban Stud.* 32, 1489–1503.
- Clay, P., 1979. Neighborhood Renewal: Middle-Class Resettlement and Incumbent Upgrading in American Neighborhoods. Lexington Books, Lexington, Massachusetts.
- Cliff, A.D., Ord, J.K., 1973. Spatial autocorrelation, Monographs in spatial and environmental systems analysis. Pion, London.
- Covington, J., Taylor, R.B., 1989. Gentrification and Crime: Robbery and Larceny Changes in Appreciating Baltimore Neighborhoods during the 1970s. *Urban Aff. Q.* 25, 142–172.
<https://doi.org/10.1177/004208168902500109>
- Curran, W., 2004. Gentrification and the Nature of Work: Exploring the Links in Williamsburg, Brooklyn. *Environ. Plan. Econ. Space* 36, 1243–1258. <https://doi.org/10.1068/a36240>
- Dacey, M., 1968. A review on measures of contiguity for two and k-color maps, in: *Spatial Analysis; a Reader in Statistical Geography*. Prentice Hall, Englewood Cliffs, NJ.
- Delmelle, E., 2016. Mapping the DNA of urban neighborhoods: Clustering longitudinal sequences of neighborhood socioeconomic change. *Ann. Am. Assoc. Geogr.* 106, 36–56.
- Desmond, M., Gromis, A., Edmonds, L., Hendrickson, J., Krywokuski, K., Leung, L., Porton, A., 2018. Eviction Lab National Database: Version 1.0.
- Ding, L., Hwang, J., 2016. The Consequences of Gentrification: A Focus on Residents' Financial Health in Philadelphia. *Cityscape* 18, 27–56.
- Ding, L., Hwang, J., Divringi, E., 2016. Gentrification and residential mobility in Philadelphia. *Reg. Sci. Urban Econ.* 61, 38–51. <https://doi.org/10.1016/j.regsciurbeco.2016.09.004>
- Duncan, M., 2008. Comparing Rail Transit Capitalization Benefits for Single-Family and Condominium Units in San Diego, California. *Transp. Res. Rec.* 2067, 120–130. <https://doi.org/10.3141/2067-14>
- Ellen, I., Ding, L., 2016. Gentrification: Advancing Our Understanding of Gentrification. *Cityscape J. Policy Dev. Res.* 18.
- Ellen, I., O'Regan, K., 2010. How Low Income Neighborhoods Change: Entry, Exit and Enhancement (Working Paper). NYU Wagner School and Furman Center for Real Estate & Urban Policy.
- Enterprise Community Partners, 2019. Gentrification Comparison Tool [WWW Document]. URL <https://www.enterprisecommunity.org/policy-and-advocacy/policy-development-and-research/gentrification-comparison-tool> (accessed 6.7.19).
- Florida, R., 2014. The Rise of the Creative Class--Revisited: Revised and Expanded. Basic Books.
- Florida, R., 2003. Cities and the Creative Class. *City Community* 2, 3–19. <https://doi.org/10.1111/1540-6040.00034>
- Florida, R., Cohen, W., 1999. Industrializing Knowledge, in: *Industrializing Knowledge: University-Industry Linkages in Japan and the United States*.
- Formoso, D., N Weber, R., S Atkins, M., 2010. Gentrification and urban children's well-being: tipping the scales from problems to promise. *Am. J. Community Psychol.* 46, 395–412.
<https://doi.org/10.1007/s10464-010-9348-3>

- Fraser, J., 2017. Housing Policy and Inclusionary Zoning Feasibility Study. Metropolitan Planning Department.
- Freeman, L., 2005. Displacement or Succession? Urban Aff. Rev. 40.
- Galster, G., Peacock, S., 1986. Urban Gentrification: Evaluating Alternative Indicators. Soc. Indic. Res. 18, 321–337.
- Gamper-Rabindran, S., Mastromonaco, R., Timmins, C., 2011. Valuing the Benefits of Superfund Site Remediation: Three Approaches to Measuring Localized Externalities (Working Paper No. 16655). National Bureau of Economic Research. <https://doi.org/10.3386/w16655>
- Gamper-Rabindran, S., Timmins, C., 2011. Hazardous Waste Cleanup, Neighborhood Gentrification, and Environmental Justice: Evidence from Restricted Access Census Block Data. Am. Econ. Rev. 101, 620–624.
- gentrification | Definition of gentrification in English by Oxford Dictionaries [WWW Document], 2019. . Oxf. Dictionaries Engl. URL <https://en.oxforddictionaries.com/definition/gentrification> (accessed 6.3.19).
- GeoLytics, Inc., 2005. CensusCD Neighborhood Change Database (NCDB): 1970-2000 Tract Data. Geolytics, East Brunswick, NJ.
- Gibbons, J., Barton, M.S., 2016. The Association of Minority Self-Rated Health with Black versus White Gentrification. J. Urban Health Bull. N. Y. Acad. Med. 93, 909–922. <https://doi.org/10.1007/s11524-016-0087-0>
- Gibbons, S., Machin, S., 2005. Valuing rail access using transport innovations. J. Urban Econ. 57, 148–169. <https://doi.org/10.1016/j.jue.2004.10.002>
- Glass, R., 1964. London: Aspects of Change. MacGibbon & Kee, London.
- Goetz, E.G., Lewis, B., Damiano, A., Calhoun, M., 2019. THE DIVERSITY OF GENTRIFICATION: The Regents of the University of Minnesota.
- Grier, G., Grier, E., 1978. Urban Displacement: A Reconnaissance. Grier Partnership, Bethesda, MD.
- Halle, D., Tiso, E., 2014. New York's New Edge: Contemporary Art, the High Line, and Urban Megaprojects on the Far West Side. University of Chicago Press, Chicago.
- Hammel, D.J., Wyly, E.K., 1996. A Model for Identifying Gentrified Areas with Census Data. Urban Geogr. 17, 248–268. <https://doi.org/10.2747/0272-3638.17.3.248>
- Hamnett, C., 1991. The Blind Men and the Elephant: The Explanation of Gentrification. Trans. Inst. Br. Geogr. 16, 173–189. <https://doi.org/10.2307/622612>
- Hampson, R., 2005. Gentrification a Boost to Everyone. USA Today.
- Harper, G., Cotton, C., 2015. Nashville Music Industry: Impact, Contribution, and Cluster Analysis. Nashville, TN.
- Haruch, S., 2014. High Rises vs. Honky Tonks. N. Y. Times.
- Helms, A.C., 2003. Understanding gentrification: an empirical analysis of the determinants of urban housing renovation. J. Urban Econ. 54, 474–498. [https://doi.org/10.1016/S0094-1190\(03\)00081-0](https://doi.org/10.1016/S0094-1190(03)00081-0)
- Hess, D.B., Almeida, T.M., 2007. Impact of Proximity to Light Rail Rapid Transit on Station-area Property Values in Buffalo, New York. Urban Stud. 44, 1041–1068. <https://doi.org/10.1080/00420980701256005>
- Housing Nashville: Nashville & Davidson County's Housing Report, 2017. . Office of the Mayor Megan Barry.
- Huynh, M., Maroko, A.R., 2014. Gentrification and preterm birth in New York City, 2008–2010. J. Urban Health Bull. N. Y. Acad. Med. 91, 211–220. <https://doi.org/10.1007/s11524-013-9823-x>
- Hwang, J., Sampson, R.J., 2014. Divergent Pathways of Gentrification: Racial Inequality and the Social Order of Renewal in Chicago Neighborhoods. Am. Sociol. Rev. 79, 726–751. <https://doi.org/10.1177/0003122414535774>

- Joint Center for Housing Studies, 2018. The State of the Nation's Housing. Harvard University.
- Kennedy, M., Leonard, P., 2001. Dealing with Neighborhood Change: A Primer on Gentrification and Policy Choices. The Brookings Institution Center on Urban and Metropolitan Policy.
- Kerner J.R., O., 1968. Report of the National Advisory Commission on Civil Disorders. Washington D.c.
- Kiviat, B., 2008. Gentrification: Not Ousting the Poor? Time.
- Knaap, G.J., Ding, C., Hopkins, L.D., 2001. Do Plans Matter?: The Effects of Light Rail Plans on Land Values in Station Areas. *J. Plan. Educ. Res.* 21, 32–39. <https://doi.org/10.1177/0739456X0102100103>
- Kreager, D.A., Lyons, C.J., Hays, Z.R., 2011. Urban Revitalization and Seattle Crime, 1982–2000. *Soc. Probl.* 58, 615–639. <https://doi.org/10.1525/sp.2011.58.4.615>
- Kuhn, M., 2019. caret.
- Landis, J.R., Koch, G.G., 1977. The Measurement of Observer Agreement for Categorical Data. *Biometrics* 33, 159–174. <https://doi.org/10.2307/2529310>
- Lang, M., 1986. Measuring Economic Benefits from Gentrification. *J. Urban Aff.* 8, 27–39. <https://doi.org/10.1111/j.1467-9906.1986.tb00152.x>
- Lang, M.H., 1982. Gentrification amid urban decline: strategies for America's older cities. Ballinger Pub. Co.
- Larsson, N., 2017. The Final Bar? How Gentrification Threatens America's Music Cities. *The Guardian*.
- Laska, S., Spain, D., 1980. Back to the City. Elsevier.
- Lee, Y.Y., 2010. Gentrification and Crime: Identification Using the 1994 Northridge Earthquake in Los Angeles. *J. Urban Aff.* 32, 549–577. <https://doi.org/10.1111/j.1467-9906.2010.00506.x>
- Lees, L., 2007. Progress in Gentrification Research? *Environ. Plan. Econ. Space* 39, 228–234. <https://doi.org/10.1068/a39329>
- Lees, L., 2003. Super-gentrification: The Case of Brooklyn Heights, New York City. *Urban Stud.* 40, 2487–2509. <https://doi.org/10.1080/0042098032000136174>
- Lees, L., 2000. A Re-appraisal of Gentrification: Towards a Geography of Gentrification. *Prog. Hum. Geogr.* 24, 389–408. <https://doi.org/10.1191/030913200701540483>
- Lester, T.W., Hartley, D.A., 2013. The Long Term Employment Impacts of Gentrification in the 1990 s March 30 th , 2013 DRAFT DO NOT CIRCULATE.
- Ley, D., 1986. Alternative Explanations for Inner-City Gentrification: A Canadian Assessment. *Ann. Assoc. Am. Geogr.* 76, 521–535. <https://doi.org/10.1111/j.1467-8306.1986.tb00134.x>
- Ley, D., 1980. Liberal Ideology and the Postindustrial City*. *Ann. Assoc. Am. Geogr.* 70, 238–258. <https://doi.org/10.1111/j.1467-8306.1980.tb01310.x>
- Liaw, A., Wiener, M., 2018. randomForest: Breiman and Cutler's Random Forests for Classification and Regression.
- Ling, C., Delmelle, E., 2016. Classifying multidimensional trajectories of neighbourhood change: a self-organizing map and k-means approach: *Annals of GIS: Vol 22, No 3* 22, 173–186. <https://doi.org/10.1080/19475683.2016.1191545>
- Maciag, M., 2015. Gentrification in America Report.
- MacQueen, J., 1967. Some methods for classification and analysis of multivariate observations. Presented at the Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics, The Regents of the University of California.
- Maloutas, T., 2012a. Contextual Diversity in Gentrification Research. *Crit. Sociol.* 38, 33–48. <https://doi.org/10.1177/0896920510380950>
- Maloutas, T., 2012b. Contextual Diversity in Gentrification Research. *Crit. Sociol.* 38, 33–48. <https://doi.org/10.1177/0896920510380950>
- Mani, A., Mullainathan, S., Shafir, E., Zhao, J., 2013. Poverty Impedes Cognitive Function. *Science* 341, 976–980. <https://doi.org/10.1126/science.1238041>

- Marcuse, P., Smith, N., Williams, P., 1986. Abandonment, Gentrification, and Displacement: The Linkages in New York City, in: *Gentrification of the City*. Routledge, New York, pp. 153–177.
- Martin, I.W., Beck, K., 2018. Gentrification, Property Tax Limitation, and Displacement. *Urban Aff. Rev.* 54, 33–73. <https://doi.org/10.1177/1078087416666959>
- McDonald, J.F., Osuji, C.I., 1995. The effect of anticipated transportation improvement on residential land values. *Reg. Sci. Urban Econ.* 25, 261–278. [https://doi.org/10.1016/0166-0462\(94\)02085-U](https://doi.org/10.1016/0166-0462(94)02085-U)
- McDonald, S.C., 1986. Does Gentrification Affect Crime Rates? *Crime Justice* 8, 163–201.
- McEwen, B.S., 1998. Stress, adaptation, and disease: Allostasis and allostatic load, in: *Molecular Aspects, Integrative Systems, and Clinical Advances*, Annals of the New York Academy of Sciences. New York Academy of Sciences, New York, NY, US, pp. 33–44.
- McGee, J., 2018. Big financial companies increasingly choosing Nashville. *Tennessean*.
- McKenzie, R.D., Park, R., Burgess, E., 1925. The Ecological Approach to the Study of the Human Community, in: *The City*. The University of Chicago Press, Chicago, pp. 63–79.
- McKinnish, T., Walsh, R., White, K., 2008. Who Gentrifies Low-Income Neighborhoods? (Working Paper No. 14036). National Bureau of Economic Research. <https://doi.org/10.3386/w14036>
- Mellor, A., Haywood, A., Stone, C., Jones, S.D., 2013. The Performance of Random Forests in an Operational Setting for Large Area Sclerophyll Forest Classification. *Remote Sens.* 5, 2838–2856. <https://doi.org/10.3390/rs5062838>
- Meltzer, R., Ghorbani, P., 2017. Does gentrification increase employment opportunities in low-income neighborhoods? *Reg. Sci. Urban Econ.* 66, 52–73. <https://doi.org/10.1016/j.regsciurbeco.2017.06.002>
- Metro Codes Department, 2019. Nashville Open Data Portal.
- Metro Codes Department, 2018. Building Permits Issued. Metro Codes Department.
- Metro Government of Nashville & Davidson County - Parks and Recreation, 2016. Parks- Metro Parks Boundary Outlines (GIS).
- Monroe Sullivan, D., Shaw, S.C., 2011. Retail Gentrification and Race: The Case of Alberta Street in Portland, Oregon. *Urban Aff. Rev.* 47, 413–432. <https://doi.org/10.1177/1078087410393472>
- Munneke, H.J., 1996. Redevelopment Decisions for Commercial and Industrial Properties. *J. Urban Econ.* 39, 229–253. <https://doi.org/10.1006/juec.1996.0013>
- Munneke, H.J., Womack, K.S., 2013. Gentrification and the decision to renovate or teardown. *Nashville Region's 2018 Vital Signs*, 2018. . Greater Nashville Regional Council.
- Nashville Tourism & Hospitality, 2018. . Nashville Convention & Visitors Corporation, Nashville, TN.
- Newman, K., Wyly, E.K., 2006. The Right to Stay Put, Revisited: Gentrification and Resistance to Displacement in New York City. *Urban Stud.* 43, 23–57. <https://doi.org/10.1080/00420980500388710>
- Phillips, D., Flores Jr, L., Henderson, J., 2015. Development without Displacement: Resisting Gentrification in the Bay Area. *Causa Justa :: Just Cause*.
- Plazas, D., 2017. The costs of growth and change in Nashville. *The Tennessean*.
- Podagrosi, A., Vojnovic, I., Pigozzi, B., 2011. The diversity of gentrification in Houston's urban renaissance: From cleansing the urban poor to supergentrification. *ResearchGate* 43.
- Pollack, S., Bluestone, B., Billingham, C., 2010. Maintaining diversity in America's transit-rich neighborhoods: Tools for equitable neighborhood change, New England Community Developments. Federal Reserve Bank of Boston.
- Reicher, M., 2017. Is Wall Street upending Nashville neighborhoods? These professors want to know. *The Tennessean*.
- Royall, E.B., 2016. Towards an epidemiology of gentrification : modeling urban change as a probabilistic process using k-means clustering and Markov models (Thesis). Massachusetts Institute of Technology.

- Rutheiser, C., 2011. The Promise and Prospects of Anchor Institutions: some thoughts on an emerging field [WWW Document]. URL https://www.huduser.gov/portal/pdredge/pdr_edge_hudpartprt_062211.html (accessed 6.23.19).
- Shaw, K., 2005. Local limits to Gentrification: implications for a new urban policy. Routledge.
- Smith, N., 1996. The New Urban Frontier: Gentrification and the Revanchist City. Psychology Press.
- Smith, N., 1987. Gentrification and the Rent Gap. *Nnals Assoc. Am. Geogr.* 77, 462–465.
- Smith, N., 1982. Gentrification and Uneven Development. *Econ. Geogr.* 58, 139–155. <https://doi.org/10.2307/143793>
- Smith, N., 1979. Toward a Theory of Gentrification A Back to the City Movement by Capital, not People. *J. Am. Plann. Assoc.* 45, 538–548. <https://doi.org/10.1080/01944367908977002>
- Stern, M.J., Seifert, S.C., 2007. Culture and Urban Revitalization: A Harvest Document. University of Pennsylvania.
- Sugar, C.A., James, G.M., 2003. Finding the Number of Clusters in a Dataset: An Information-Theoretic Approach. *J. Am. Stat. Assoc.* 98, 750–763.
- Tatian, P.A., Kingsley, G.T., Parilla, J., Pendall, R., 2016. Building Successful Neighborhoods [WWW Document]. Urban Inst. URL <https://www.urban.org/research/publication/building-successful-neighborhoods> (accessed 6.6.19).
- Tobler, W.R., 1970. A Computer Movie Simulating Urban Growth in the Detroit Region. *Econ. Geogr.* 46, 234. <https://doi.org/10.2307/143141>
- U. S. Census Bureau, 2017. 2017 TIGER/Line Shapefiles.
- United Nations Department of Economic and Social Affairs, 2018. 68% of the world population projected to live in urban areas by 2050, says UN | UN DESA | United Nations Department of Economic and Social Affairs [WWW Document]. URL <https://www.un.org/development/desa/en/news/population/2018-revision-of-world-urbanization-prospects.html> (accessed 12.19.18).
- United States Center for Disease Control, 2017. CDC - Healthy Places - Health Effects of Gentrification [WWW Document]. URL <https://www.cdc.gov/healthyplaces/healthtopics/gentrification.htm> (accessed 12.19.18).
- United States Housing and Urban Development, 2018. Displacement of Lower-Income Families in Urban Areas Report 26.
- U.S. Census Bureau, 2018. Glossary [WWW Document]. URL <https://www.census.gov/programs-surveys/geography/about/glossary.html> (accessed 6.22.19).
- U.S. Census Bureau, 2016. American Community Survey 2012 - 2016.
- US EPA, 2017. Population Surrounding 1,836 Superfund Remedial Sites [WWW Document]. URL <https://www.epa.gov/sites/production/files/2015-09/documents/webpopulationrsuperfundsites9.28.15.pdf>
- Vicario, L., Martinez Monje, P.M., 2003. Another “Guggenheim Effect”? The Generation of a Potentially Gentrifiable Neighbourhood in Bilbao. *Urban Stud.* 40, 2383–2400. <https://doi.org/10.1080/0042098032000136129>
- Viera, A.J., Garrett, J.M., 2005. Understanding interobserver agreement: the kappa statistic. *Fam. Med.* 37, 360–363.
- Vigdor, J.L., Massey, D.S., Rivlin, A.M., 2002. Does Gentrification Harm the Poor? [with Comments]. *Brook.-Whart. Pap. Urban Aff.* 133–182.
- Voith, R., 1993. Changing Capitalization of CBD-Oriented Transportation Systems: Evidence from Philadelphia, 1970–1988. *J. Urban Econ.* 33, 361–376. <https://doi.org/10.1006/juec.1993.1021>
- Walker, K., 2019. tidy census.
- Walker, K., Rudis, B., 2019. tigris.

- Ward, G., Reicher, M., 2017. Nashville neighborhoods get new names amid gentrification. *The Tennessean*.
- Wardrip, K., 2011. Public Transit's Impact on Housing Costs: A Review of the Literature, Insights from Housing Policy Research. Center For Housing Policy.
- Weber, R., Doussard, M., Bhatta, S.D., Mcgrath, D., 2006. Tearing the City Down: Understanding Demolition Activity in Gentrifying Neighborhoods. *J. Urban Aff.* 28, 19–41.
<https://doi.org/10.1111/j.0735-2166.2006.00257.x>
- Weissbourd, R., Bodini, R., He, M., 2009. Dynamic Neighborhoods: New Tools for Community and Economic Development. *Living Cities: The National Community Development Initiative*.
- White, G.B., 2015. The Many Meanings of Gentrification [WWW Document]. *The Atlantic*. URL <https://www.theatlantic.com/business/archive/2015/05/the-g-word-gentrification-and-its-many-meanings/394016/> (accessed 6.7.19).
- Wickham, H., Francois, R., Henry, L., Muller, K., 2019. dplyr: A Grammar of Data Manipulation.
- Wilsem, J.A. van, Wittebrood, K.A., Graaf, N.D. de, 2006. Socioeconomic Dynamics of Neighborhoods and the Risk of Crime Victimization: A Multilevel Study of Improving, Declining, and Stable Areas in the Netherlands. 247.
- Wood, D., Halfon, N., Scarlata, D., Newacheck, P., Nessim, S., 1993. Impact of family relocation on children's growth, development, school function, and behavior. *JAMA* 270, 1334–1338.
- Wyly, E., Newman, K., Schafran, A., Lee, E., 2010. Displacing New York. *Environ. Plan. Econ. Space* 42, 2602–2623. <https://doi.org/10.1068/a42519>
- Ye, M., Vojnovic, I., Chen, G., 2015. The landscape of gentrification: exploring the diversity of “upgrading” processes in Hong Kong, 1986–2006. *Urban Geogr.* 36, 471–503.
<https://doi.org/10.1080/02723638.2015.1010795>
- Yi, B., Qiao, H., Yang, F., Xu, C., 2010. An Improved Initialization Center Algorithm for K-Means Clustering, in: 2010 International Conference on Computational Intelligence and Software Engineering. Presented at the 2010 International Conference on Computational Intelligence and Software Engineering, pp. 1–4. <https://doi.org/10.1109/CISE.2010.5676975>
- Zuk, M., Bierbaum, A., Chapple, K., Gorska, K., Loukaitou-Sideris, A., Ong, P., Thomas, T., 2015. Gentrification, Displacement and the Role of Public Investment: A Literature Review. University of California.

Appendix

In addition to the random forest model featured in the main sections, we also conducted a principle component analysis and logistic regression model to differentiate gentrified tracts from non-gentrified tracts. This process is discussed below.

Preprocessing

Principal Component Analysis (PCA) is a statistical procedure for reducing dimensionality while retaining maximum information via orthogonal transformations of correlated and uncorrelated variables called principal components. PCA identifies a linear combination of variables that explain the maximum variance of the dataset by projecting data onto lower dimensions. The first principle component minimizes the distance between original data and the newly projected axis and simultaneously maximizes the variance between the projected points. Subsequent principle components are processed likewise, each attempting to explain a portion of the remaining variance in the dataset (Lever et al. 2017). The resulting principle components are orthogonal and therefore uncorrelated linear combinations of the original data. We use PCA in this study to accommodate the large number of predictive variables and account for correlation among them within our logistic regression model.

PCA is necessary measure to avoid the “curse of dimensionality”, a frequent problem encountered in modeling high-dimensional datasets. This phenomena describes model overfitting and a loss in predictive performance as more input variables or dimensions are added. As the dimensionality of the dataset grow, the number of observations used in training and testing the data represent a smaller percentage of the feature space. Machine learning classifiers attempt to draw distinct boundaries and create specific exceptions for each observation within the training data in high dimensionality spaces. The problem arises when the model is applied to real world or test data that the model has not previously seen. The highly detailed exceptions of the trained model attempt to reconstruct the inherent noise in the training data set and fail to distinguish the underlying structure of the real world data. Lacking generalizability, the model has been over fitted by adding too many dimensions and not enough training data to justify the number of predictive variables.

Our analysis focuses on a fixed amount of census tracts (161) constituting Davidson County, Tennessee. We are not afforded the benefit of adding additional training data and as such make use of PCA to reduce the dimensionality of the dataset. We are fundamentally interested in a variety of social, economic, spatial, and built environment predictors that may have distinguish a neighborhoods outcome over time. The inclusion of the 32 selected variables was guided by previous literature, domain knowledge, and data availability.

We centered and scaled all predictors and then conducted PCA on our training dataset. Figure 25 illustrates the percent of variance explained by the first ten principle components. We are then motivated to select a subset of principle components that retain maximum information in the metric of variance. The first five principle components explain 73% of the variance in the dataset. Jolliffe (2002) advises to truncate the list of components so that the selected components $70\% < R^2 < 90\%$.

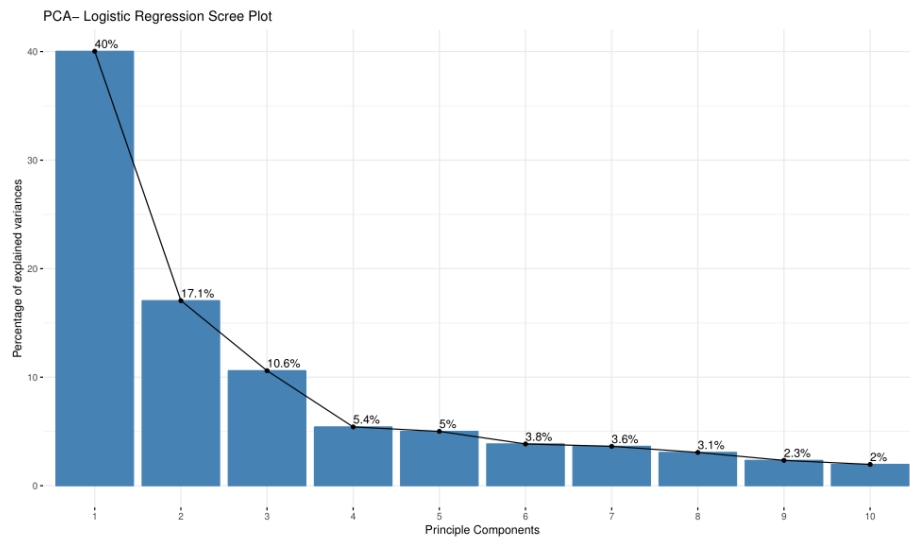


Figure 25. Principle component analysis variance scree plot

Figure 26 shows the loadings associated with each predictive variables within the first two PC dimensions. This plot graphs the coefficients of each variable for the first principle component against its coefficient within the second principle component. Positive correlation between vectors is indicated by vectors pointing in the similar direction. Vectors in opposite directions (180 degrees) suggests negative correlation while vectors of 90 degrees suggests no correlation. Vectors that extend horizontally can be understood as significant to the first principle component while more vertical vectors indicate heavier loading on the second PC.

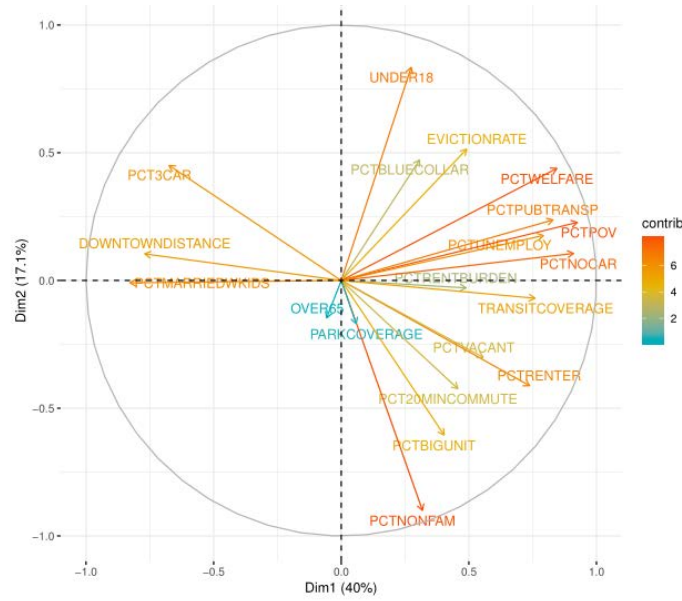


Figure 26. PCA biplot.

Logistic Regression

Binary logistic regression is one type statistical model used by previous researchers of neighborhood change (Chapple 2009; Winston and Walker 2012). Logistical regression models use input variables to develop a probability between 0 and 1 that an observation belongs to a positive class. The logistic function that gives values ranging from 0 to 1.

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}$$

$$\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}$$

The terms may be rearranged into figure x with $p(X)/[1-p(X)]$ left hand side representing the odds ratio as opposed to the probability of X. We can take the logarithm of both sides of the equation to yield the logit or log-odds function:

$$\log\left(\frac{p(X)}{1 - p(X)}\right) = \beta_0 + \beta_1 X$$

The coefficients β_0 and β_1 are estimated using a labeled training data subset such that the predicted probability $P^{\wedge}(X_i)$ follows each observation's class status. Positive classes are indicated with $p(X)$ closer to one while negative classes will ideally fall closer to zero (Gareth et al. 2013). After the coefficient

estimation, we can easily compute the probability of an observation falling into a positive or negative class (above first eq.).

Thresholds Logistical regression estimates the conditional probability that $Y = 1$ given the X variables, or $\Pr(Y = 1 | X = x)$. The logit is a link function that describes the probability of a given outcome.

Logistic Regression Model

We fit a principle components logistic regression model using five uncorrelated predictive variables in order to predict the outcome classification of “gentrifying”. The model was fit using a training dataset comprised of 60 percent of the observations with leave-one-out cross validation. Table 10 provides the fitted coefficients of the variables. PC1 and PC3 were reported as statistically significant at a significance threshold of .05.

Table 9 provides the eigenvectors associated with each principle component. Positive loadings are colored black while negative loadings are symbolized in red. Values greater than .2 or less than -.2 are symbolized with bold font. The first principle component loads heavily on economically vulnerable variables like poverty, unemployment, welfare, and transportation access. It also appears to capture household status dynamics such as percent renters and married families with children. Lastly PC1 loads on spatial associations that are potentially tracts located close to downtown and with ample transit coverage. The second principle component appears to focus primarily on homeowners with families and children under the age of 18. PC2 also loads on vehicle accessibility. Principle Component 3 may be indicative of rural renters with longer commutes.

Table 9. PCA eigenvalues.

Variables	PC1	PC2	PC3	PC4	PC5
PCTRENT	0.261	-0.223	0.325	-0.016	0.01
PCTVACANT	0.197	-0.163	0.103	0.278	-0.069
PCTNONFAM	0.112	-0.487	0.069	0.107	-0.035
UNDER18	0.097	0.452	0.117	-0.233	0.106
OVER65	-0.02	-0.079	-0.511	0.128	0.057
PCTPOV	0.327	0.123	0.045	-0.093	0.063
PCTUNEMPLOY	0.28	0.095	-0.022	-0.12	-0.09
PCTNOCAR	0.322	0.057	-0.011	-0.174	0.115
PCT3CAR	-0.238	0.244	-0.195	0.016	-0.122
PCTPUBTRANSP	0.293	0.128	-0.031	-0.185	0.082
PCT20MINCOMMUTE	0.161	-0.229	-0.4	-0.147	-0.174
PCTBLUECOLLAR	0.108	0.256	0.158	0.494	-0.146
PCTWELFARE	0.299	0.238	-0.01	-0.107	0.094
PCTBIGUNIT	0.142	-0.328	0.374	-0.145	0.027
PCTMARRIEDWKIDS	-0.293	-0.005	0.256	0.126	-0.052
PCTRENTBURDEN	0.173	-0.016	-0.171	0.404	-0.149
PARKCOVERAGE	0.021	-0.092	-0.117	0.262	0.898
DOWNTOWNDISTANCE	-0.272	0.056	0.33	-0.041	0.099
TRANSITCOVERAGE	0.268	-0.038	-0.027	0.15	-0.172
EVICTIIONRATE	0.174	0.278	0.141	0.426	-0.015

Table 10. PCA-logistic classification model parameters.

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.866	0.63	-4.54	0.0000
PC1	0.505	0.15	3.31	0.0009
PC2	-0.311	0.17	-1.81	0.0707
PC3	-0.830	0.35	-2.35	0.0189
PC4	0.211	0.32	0.66	0.5086
PC5	0.145	0.22	0.65	0.5187

We tested our logistic regression classification on the remaining 40 percent of holdout testing data. The testing dataset was comprised of 60 labeled observations. 52 tracts in the testing dataset did not gentrify according to our classification while 8 tracts fell into the positive class of “gentrifying”. Table 1 presents the confusion matrix with correct classifications falling on the diagonal upper left to lower right, while misclassifications on the opposite diagonal. A classification threshold of .77 was chosen to specify the final predicted outcome, which balanced sensitivity and specificity.

Table 11. PCA-logistic regression confusion matrix.

Predicted	Observed	
	Not Gentrifying	Gentrifying
Not Gentrifying	44	1
Gentrifying	8	7

The receiver operating characteristic or ROC curve is a visualization of type I (false positive) and type II (false negative) errors across all possible classification thresholds. The plot visualizes the classification trade-off between sensitivity and specificity- a higher threshold will identify more of the truly gentrified areas at the expense of falsely labeling non-gentrifying tracts as gentrifying. The area under the curve (AUC) can be used to assess the overall performance of a classifier. A perfect classifier will obtain an AUC value of 1 while a classifier that performs on par with a random chance would result in an AUC value of .5 (Gareth et al. 2013).

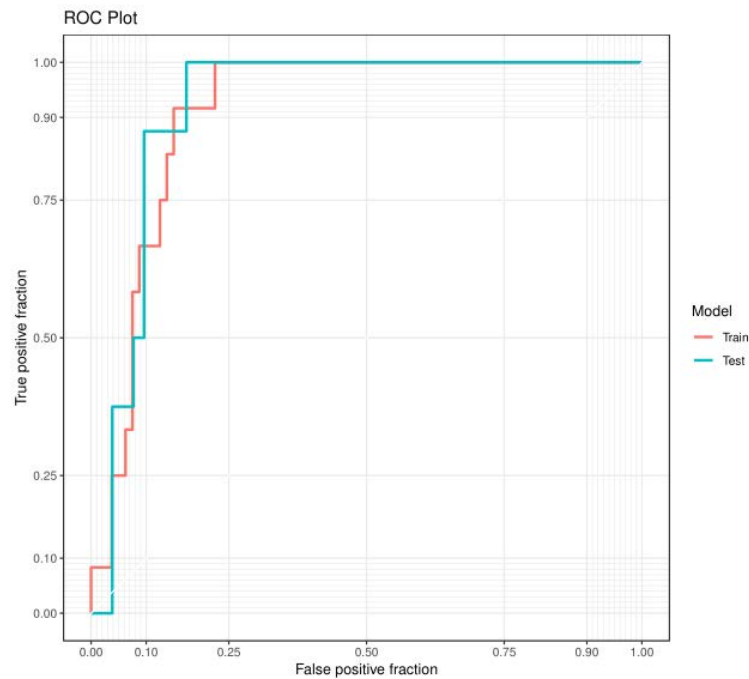


Figure 27. PCA-logistic regression receiver operator characteristic (ROC) curve.

