

# Assessing the Status of Gentrification in Metropolitan Areas in the United States

Vishwesh Srinivasan, Vanessa Venkataraman, and Hanzhen Wang

A  
G  
E  
N  
D  
A

01  
02  
03  
04  
05  
06

Introduction  
Literature Review  
Data Description  
Our Hypothesis  
Models and their Results  
Conclusions and Future Work

# Introduction

“Gentrification is the process by which central urban neighborhoods that have undergone disinvestments and economic decline experience a reversal, reinvestment, and the in-migration of a relatively well-off middle- and upper middle-class population.” - Neil Smith

Positives of Gentrification: increased investment, commercial development, increased property values, improvement in economic opportunity

Negatives of Gentrification: high rate of displacement in communities, higher housing costs, higher property taxes, forced relocation, community cultural changes, changes in businesses, dislocation

Gentrification disproportionately impacts low-income families and people of color.

In order to counteract the negative impacts of gentrification, it is important to recognize which neighborhoods are susceptible to gentrification and which communities are most vulnerable to displacement.

# Literature Review

Kolko (2007) studied the determinants of gentrification using the Neighborhood Change Database. His research supported the fact that low-income city neighborhoods closer to the city center and with older housing stock witness higher income growth.

Kolko (2009), in another study, evaluated the effect of employment location on gentrification. His findings suggested that an increase in the mean neighborhood job pay influences the tract's income growth, by which he inferred that people follow jobs.

Boustan et al. (2019) studied if condominium structure led to gentrification in central cities. They did not see a statistically significant association between the condo density and the typical characteristics of gentrified regions – income levels and education.

Rucks-Ahidiana (2021) studied the effect of racial composition on the way gentrification happens in a community. He found that the gentrifying tracts with residents of color, compared to tracts with white people, display more characteristics of poverty.

# Literature Review

Rigolon & Németh (2019) took a broad approach and considered the effect of three different categories of variables – people, place, and policy on gentrification.

They used the criteria given by Chapple et al. (2017) to identify the areas susceptible to gentrification in 2000. Then they again used the definition Chapple et al. (2017) gave to identify the areas that had actually been gentrified in 2015.

Finally, they tested how all the variables they had selected in the three categories for their model are beneficial in predicting whether a gentrification susceptible region will get gentrified or not.

They tested their model in five Combined Statistical Areas (CSAs). These CSAs are Chicago, Los Angeles, New York, San Francisco, and Washington, DC.

Their research showed that all three categories assist in predicting the probability of gentrification of the tracts. Some of the significant factors were: race, distance from downtown, units older than 30 years and population density.

# Data Description - Neighborhood Change Database (NCDB)

Census tracts are locally determined geographic units, typically including between 2,500 and 8,000 persons. Tracts are meant to approximate "neighborhoods" by capturing a group of residents with similar social characteristics, economic status, and housing conditions.

A metropolitan area or metro is a region that consists of a densely populated urban agglomeration and its surrounding territories sharing industries, commercial areas, transport network, infrastructures and housing.

A metro area usually comprises multiple principal cities, jurisdictions and municipalities: neighborhoods, townships, boroughs, cities, towns, exurbs, suburbs, counties, districts, as well as even states and nations like the euro districts. As social, economic and political institutions have changed, metropolitan areas have become key economic and political regions.

A powerful feature of the NCDB is its ability to match tracts across all the five census years, enabling users to examine changing tract characteristics between 1970 and 2010.

Since the boundaries of many tracts change between the decennial censuses (either splitting into several tracts, merging with other tracts, or appearing for the first time), a methodology has been developed to link such tracts and their associated data to standard geographic boundaries.

# Data Description - Neighborhood Change Database (NCDB)

## *Classification of Variables*

In accordance with Census Bureau conventions, NCDB variables are separated into two general groups: population and housing. Population variables are further classified into nine categories:

- General Population Characteristics
- Family Structure/Marriage
- Mobility/Transportation
- Education
- Employment/Labor Market
- Poverty/Public Assistance
- Income and Earnings
- Age Distribution
- Language Ability

While housing variables are grouped into four categories:

- Housing Tenure/Occupancy
- Housing Characteristics/Utilities
- Housing Costs/Affordability - Owners
- Housing Costs/Affordability - Renters

Additionally, there is a category for geographic identifier variables.



# Our Hypothesis



The three factors that contribute the most to determining if a neighborhood is susceptible to gentrification are...



## PEOPLE



## PLACE



## POLICY

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_{people} \cdot features_{people} + \beta_{place} \cdot features_{place} + \beta_{policy} \cdot features_{policy} \quad (1)$$

$\beta_i$  – represents the set of coefficients for each feature under the  $i$  category

$features_i$  – represents the set of features under the  $i$  category

# Our Hypothesis - Model design and Output Variable (Y)

We selected the following two regions for our analysis:

- New York-Northern New Jersey-Long Island, NY-NJ-PA Metropolitan Statistical Area (MSA)
- Los Angeles-Long Beach-Santa Ana, CA Metropolitan Statistical Area (MSA)

We choose to study the process of gentrification between the years 2000 and 2010. To do so, we use different variables from the year 2000 as predictors (X) and create the labels (Y) as 'Gentrified' (1) or 'Not Gentrified' (0).

We use the following definition to identify if a tract is gentrified or not in the year 2010:

- Change in median household income > Change in MSA median
- Change in % college educated > Change in MSA percentage
- Change in median gross rent > Change in MSA median  
(OR)
- % Increase of home value > % increase in MSA median

(Change between the year 2000 and 2010. Given by Rigolon & Németh (2019), adapted by them from Chapple et al. (2017))

# Our Hypothesis - Predictors (People)

SHRWHT0	Proportion of White population
SHRBLKO	Proportion of Black/African American population
SHRAMIO = SHRAMION/SHROD	Proportion of American Indian/Alaska Native population
SHRASNO = SHRASNOD/SHROD	Proportion of Asian population
SHRHIP0 = SHRHIPON/SHROD	Proportion of Native Hawaiian and other Pacific Islander population
SHRHSP0	Proportion of Hispanic/Latino population
ADULT0	Proportion of persons who are adults 18+ years old
CHILDO	Proportion of persons who are children under 18 years old
OLDO	Proportion of persons who are 65+ years old
MDFAMYO	Median family income last year (\$)

# Our Hypothesis - Predictors (Place)

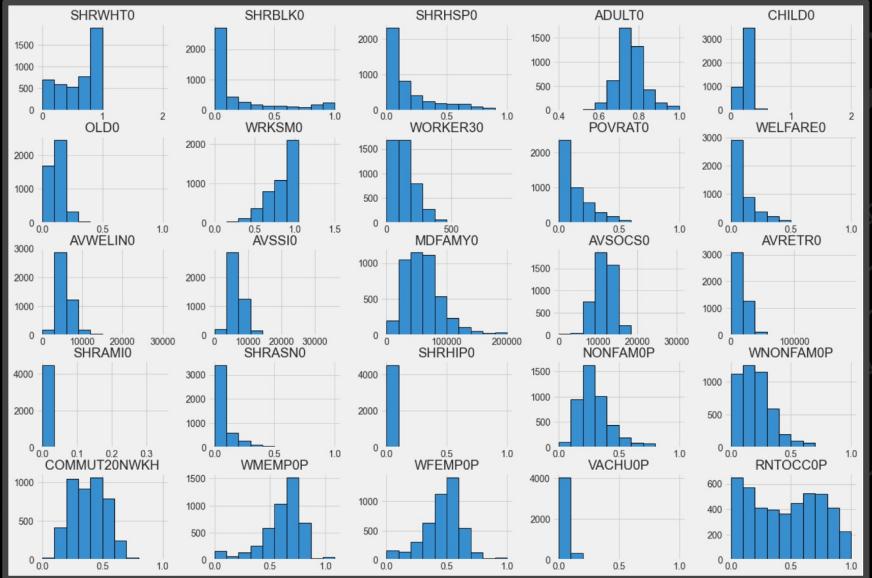
POVRATO	Proportion of total persons below the poverty level last year
WFEMPOP = WFEMPO/WF16PO	Proportion of Employed White females 16+ years old
WMEMPOP = WMEMPO/WM16PO	Proportion of Employed White males 16+ years old
WORKER30	Families with 3+ workers last year
WRKSM0	Proportion of workers who live in a metro area and work within the same metro area
NONFAM0P = NONFAM0/NUMHHS0	Proportion of Nonfamily households
WNONFAM0P = WNONFAM0/NUMHHS0	Proportion of White nonfamily households
COMMUT20NWKH = (COMMUT20 - WKHOME0)/WRCNTY0D	Proportion of Workers 16+ years old not working from home and with travel time to work less than 25 minutes
VACHUOP = VACHU0/TOTHSUN0	Proportion of vacant housing units
RNTOCC0P = RNTOCC0/TOTHSUN0	Proportion of renter-occupied housing units

# Our Hypothesis - Predictors (Policy)

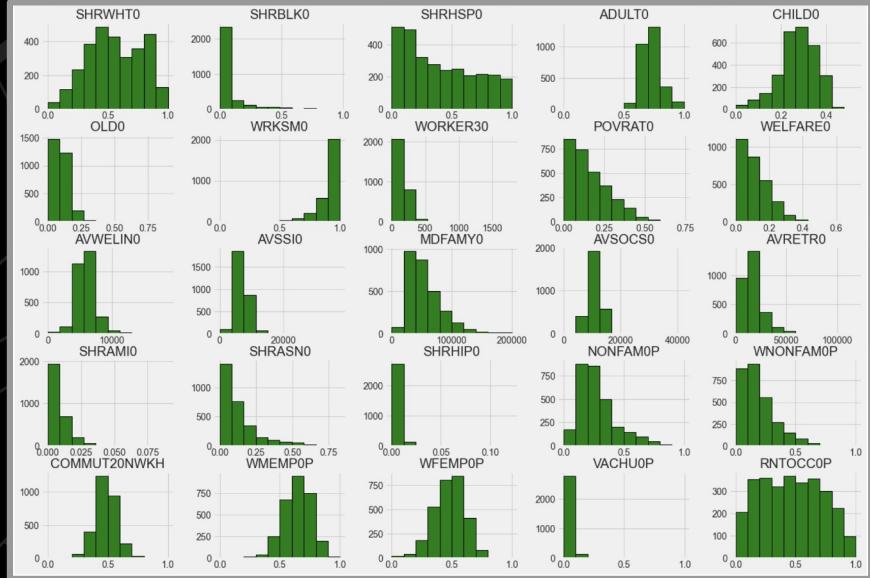
WELFARE0	Proportion of households with public assistance income (incl. SSI) last year
AVWELINO	Average public assistance income for households on public assistance last year (\$)
AVSSIO	Average Supplemental Security Income for households receiving SSI last year (\$)
AVSOC0	Average household social security income last year (\$)
AVRETRO	Average household retirement income last year (\$)

# Data Visualizations

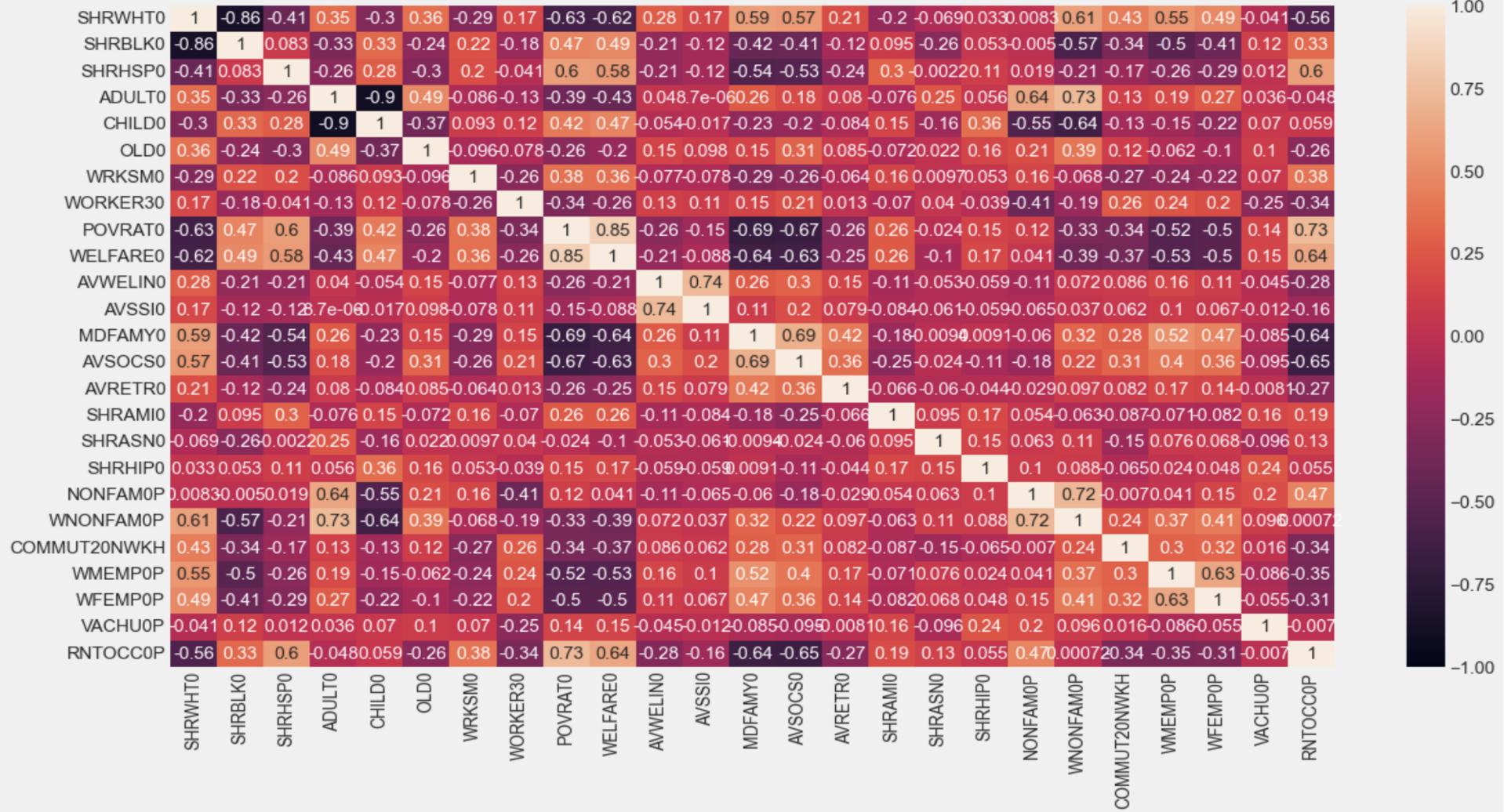
## New York Metropolitan Area

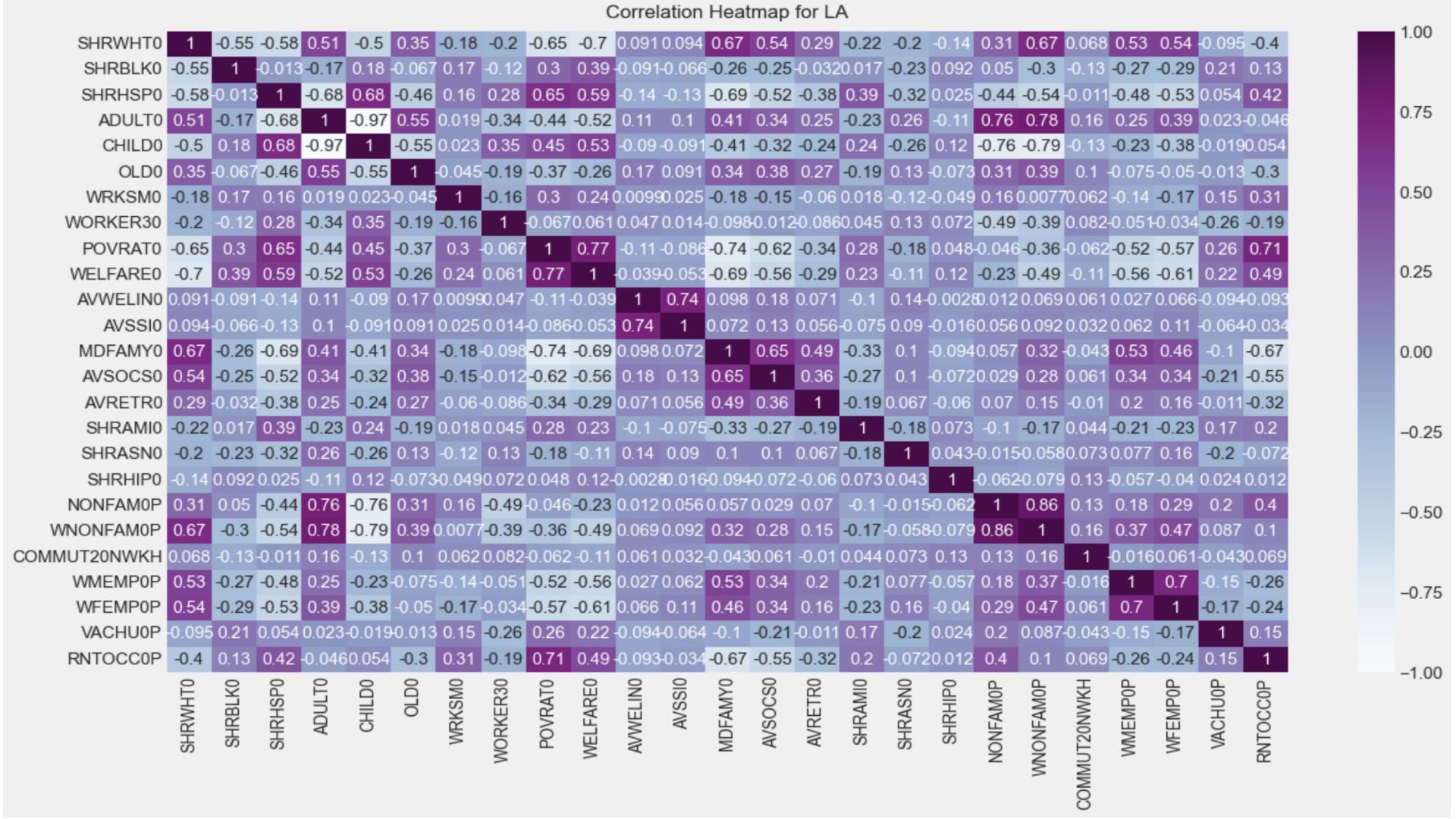


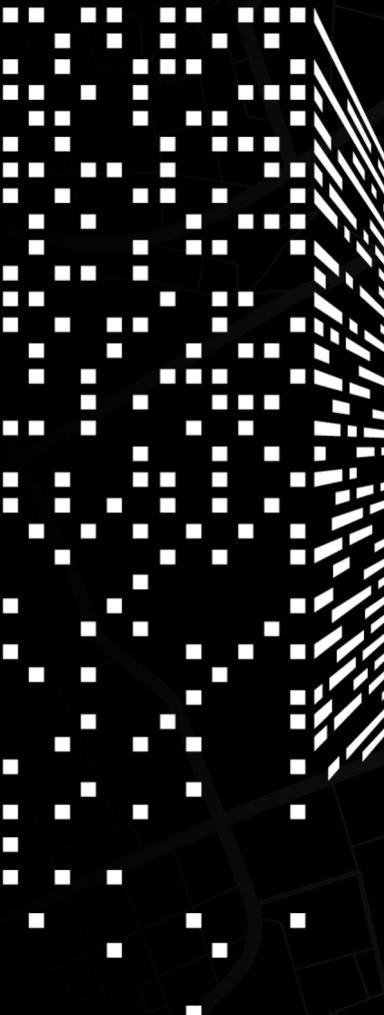
## Los Angeles Metropolitan Area



Correlation Heatmap for NY







# Models and their Results

## Data Composition:

- New York MSA: 990 tracts gentrified, 3522 tracts not gentrified
- Los Angeles MSA: 583 tracts gentrified, 2343 tracts not gentrified

## Modeling Process:

1. Normalized the variables which were not in proportions (for step 5)
2. Built a logistic regression model
3. Identified the significant variables out of all the predictors
4. Used the identified variables to build a random forest classifier
5. Used ADASYN (Adaptive Synthetic Sampling) technique to artificially generate more data of the minority class
6. Used the balanced dataset to build a random forest classifier
7. Compared the results obtained

# Models and their Results (New York MSA) - Logistic Regression Model

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-12.9445	9.2373	-1.401	0.16112
SHRWHT0	-1.4586	1.2977	-1.124	0.26099
SHRBLOK0	-1.3948	1.3207	-1.056	0.29092
SHRHSP0	-0.1774	0.7458	-0.238	0.81194
ADULT0	14.6925	9.6997	1.515	0.12984
CHILDO	10.8847	9.6264	1.131	0.25817
OLDO	-4.3342	0.8668	-5.000	5.73e-07 ***
WRKSM0	0.8315	0.2840	2.928	0.00342 **
WORKER30	0.3824	0.5380	0.711	0.47724
POVRATO	0.4840	0.8481	0.571	0.56822
WELFARE0	-0.8786	0.9698	-0.906	0.36492
AVWELINO	0.5912	0.9110	0.649	0.51636
AVSSIO	-0.5125	0.9145	-0.560	0.57518
MDFAMY0	-2.8647	0.5195	-5.514	3.50e-08 ***
AVSOC0	0.8741	0.8578	1.019	0.30820
AVRETRO	0.5164	0.7882	0.655	0.51233
SHRAMIO	-5.1414	6.3257	-0.813	0.41634
SHRASNO	-1.7052	1.3919	-1.225	0.22052
SHRHPO	-32.6760	16.0583	-2.035	0.04187 *
NONFAMOP	2.1535	0.9342	2.305	0.02116 *
WNONFAMOP	1.5820	0.9462	1.672	0.09455 .
COMMUT20NWKH	-1.7133	0.3605	-4.752	2.01e-06 ***
WMEMPOP	0.1373	0.3554	0.386	0.69933
WFEMPOP	0.8382	0.3965	2.114	0.03454 *
VACHUOP	1.1432	0.5732	1.994	0.04612 *
RNTOCCOP	-2.8206	0.3834	-7.357	1.88e-13 ***
---				

Confusion Matrix and Statistics

		Reference	
		Prediction	0 1
0	692	184	
1	14	13	

Accuracy : 0.7807

95% CI : (0.7523, 0.8073)

No Information Rate : 0.7818

P-Value [Acc > NIR] : 0.551

Kappa : 0.067

McNemar's Test P-Value : <2e-16

Sensitivity : 0.06599

Specificity : 0.98017

Pos Pred Value : 0.48148

Neg Pred Value : 0.78995

Precision : 0.48148

Recall : 0.06599

F1 : 0.11607

Prevalence : 0.21816

Detection Rate : 0.01440

Detection Prevalence : 0.02990

Balanced Accuracy : 0.52308

'Positive' Class : 1

# Models and their Results (New York MSA) - Random Forest Classifier

```
[1] "Before Data Resampling"  
Confusion Matrix and Statistics
```

Reference		
Prediction	0	1
0	681	177
1	25	20

Accuracy : 0.7763

95% CI : (0.7477, 0.8031)

No Information Rate : 0.7818

P-Value [Acc > NIR] : 0.6734

Kappa : 0.0916

McNemar's Test P-Value : <2e-16

Sensitivity : 0.10152

Specificity : 0.96459

Pos Pred Value : 0.44444

Neg Pred Value : 0.79371

Precision : 0.44444

Recall : 0.10152

F1 : 0.16529

Prevalence : 0.21816

Detection Rate : 0.02215

Detection Prevalence : 0.04983

Balanced Accuracy : 0.53306

'Positive' Class : 1

```
[1] "After Data Resampling"  
Confusion Matrix and Statistics
```

Reference		
Prediction	0	1
0	570	119
1	136	78

Accuracy : 0.7176

95% CI : (0.687, 0.7468)

No Information Rate : 0.7818

P-Value [Acc > NIR] : 1.0000

Kappa : 0.1972

McNemar's Test P-Value : 0.3164

Sensitivity : 0.39594

Specificity : 0.80737

Pos Pred Value : 0.36449

Neg Pred Value : 0.82729

Precision : 0.36449

Recall : 0.39594

F1 : 0.37956

Prevalence : 0.21816

Detection Rate : 0.08638

Detection Prevalence : 0.23699

Balanced Accuracy : 0.60165

'Positive' Class : 1

# Models and their Results (Los Angeles MSA) - Logistic Regression Model

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	1.78803	4.22208	0.423	0.6719
SHRNWTO	3.85453	1.90171	2.027	0.0427 *
SHRWHTO	3.08880	1.82588	1.692	0.0907 .
SHRHSP0	0.95846	1.07062	0.895	0.3707
ADULT0	-4.81271	4.46384	-1.078	0.2810
CHILDO	-9.46631	4.53233	-2.089	0.0367 *
OLDO	-1.99893	1.37455	-1.454	0.1459
WRKSM0	1.46492	0.67769	2.162	0.0306 *
WORKER30	-0.34712	1.47142	-0.236	0.8135
POVRAUTO	-6.69963	1.43379	-4.673	2.97e-06 ***
WELFARE0	-0.35601	1.51817	-0.235	0.8146
AVWELINO	-1.06456	0.96546	-1.103	0.2702
AVSSIO	-0.11184	1.36278	-0.082	0.9346
MDFAMYO	-4.31326	0.95785	-4.503	6.70e-06 ***
AVSOC50	-0.18423	1.43471	-0.128	0.8978
AVRETRO	-1.28245	0.84323	-1.521	0.1283
NONFAMOP	-0.96935	1.45259	-0.667	0.5046
WNONFAMOP	1.52376	1.55486	0.980	0.3271
COMMUT20NWKH	0.71423	0.64205	1.112	0.2660
WMEMPOP	1.22875	0.80505	1.526	0.1269
WFEMPOP	0.66118	0.80576	0.821	0.4119
VACHUOP	-1.00025	1.52945	-0.654	0.5131
RNTOCCOP	-0.05735	0.50697	-0.113	0.9099
---				

Confusion Matrix and Statistics

		Reference	
		Prediction	0 1
0	454	116	
1	11	5	

Accuracy : 0.7833  
 95% CI : (0.7477, 0.816)  
 No Information Rate : 0.7935  
 P-Value [Acc > NIR] : 0.7482

Kappa : 0.026  
 Mcnemar's Test P-Value : <2e-16

Sensitivity : 0.041322  
 Specificity : 0.976344  
 Pos Pred Value : 0.312500  
 Neg Pred Value : 0.796491  
 Precision : 0.312500  
 Recall : 0.041322  
 F1 : 0.072993  
 Prevalence : 0.206485  
 Detection Rate : 0.008532  
 Detection Prevalence : 0.027304  
 Balanced Accuracy : 0.508833

'Positive' Class : 1

# Models and their Results (Los Angeles MSA) - Random Forest Classifier

```
[1] "Before Data Resampling"  
Confusion Matrix and Statistics
```

Reference		
Prediction	0	1
0	435	109
1	30	12

Accuracy : 0.7628

95% CI : (0.7262, 0.7967)

No Information Rate : 0.7935

P-Value [Acc > NIR] : 0.9689

Kappa : 0.0457

Mcnemar's Test P-Value : 3.694e-11

Sensitivity : 0.09917

Specificity : 0.93548

Pos Pred Value : 0.28571

Neg Pred Value : 0.79963

Precision : 0.28571

Recall : 0.09917

F1 : 0.14724

Prevalence : 0.20648

Detection Rate : 0.02048

Detection Prevalence : 0.07167

Balanced Accuracy : 0.51733

'Positive' Class : 1

```
[1] "After Data Resampling"  
Confusion Matrix and Statistics
```

Reference		
Prediction	0	1
0	367	66
1	98	55

Accuracy : 0.7201

95% CI : (0.6819, 0.7562)

No Information Rate : 0.7935

P-Value [Acc > NIR] : 0.99999

Kappa : 0.2221

Mcnemar's Test P-Value : 0.01549

Sensitivity : 0.45455

Specificity : 0.78925

Pos Pred Value : 0.35948

Neg Pred Value : 0.84758

Precision : 0.35948

Recall : 0.45455

F1 : 0.40146

Prevalence : 0.20648

Detection Rate : 0.09386

Detection Prevalence : 0.26109

Balanced Accuracy : 0.62190

'Positive' Class : 1

# Conclusions

From the '**People**' category, we find the following variables statistically significant:

- We find '**race**' to be a significant variable, in NY MSA, just the proportion of Native Hawaiian and other Pacific Islander population, whereas in LA MSA, along with this proportion of population, the proportion of Black, American Indians and Asians are also significant.
- In NY MSA, we find the proportion of 65+ years significant whereas in LA MSA, we find the proportion of children under 18+ to be significant.
- We find the '**median family income**' to be significant for both the NY and LA MSAs.

From the '**Place**' category, we find the following variables statistically significant:

- We find the proportion of people living in the same metro area they are working in to be significant in both the NY and LA MSAs. In addition in NY MSA, we find that the proportion of people who travel less than 25 minutes daily to work also significant.

# Conclusions

- In LA MSA, the proportion of people under the poverty level also to be significant.
- In NY MSA, variables like the proportion of Nonfamily households, proportion of employed white females, proportion of vacant and rented households are significant.

From the '**Policy**' category, we do not find any of the variables that we considered to be statistically significant.

# Future Work

- To study the effect of other variables which we did not consider from the NCDB. These new variables could act as strong predictors and improve the classification model's performance.
- We have very limited variables pertaining to 'Policy' in the NCDB, hence we can explore other datasets which can give information about development programs and study their effect on the process of gentrification.
- Instead of having a binary output variable - whether a region will get gentrified or not, we could possibly change the output variable to a 'gentrification score,' which represents the degree of gentrification.



# Q & A



Thank you!