

IDENTIFYING METROPOLITAN AREAS IN THE SOUTH EASTERN UNITED STATES SIMILAR TO SILICON VALLEY

TIMOTHY BRANNON

MARCH 2019

INTRODUCTION

An upstart Information Technology company has grown rapidly and is looking to move into an office space that will serve as its Headquarters. The company is based in the South East and would like to identify metropolitan areas in the South Eastern United States as possible location candidates. The ideal location would be similar to Silicon Valley.

Data Question:

Which metropolitan areas in North Carolina, Virginia, Georgia, South Carolina and Eastern Tennessee are most comparable to Silicon Valley?

SILICON VALLEY TARGET

Target City: San Jose, Ca

Control Cities: Palo Alto & Santa Clara

3.

SILICON VALLEY TARGET

PALO ALTO

67 THOUSAND

SANTA CLARA

126 THOUSAND

SAN JOSE

1.04 MILLION

Silicon Valley consists of several cities in California, with San Jose being the most populated. This was chosen as the target city of comparison. Santa Clara and Palo Alto were included as controls.

DATA STRATEGY

STEP ONE

Identify
Prospective Cities

STEP TWO

Compile
Demographic
data

STEP FIVE

Scale & Preprocess
Feature Set

STEP FOUR

Analyze, Review
& Feature Select

STEP THREE

Compile Local
Venue data

STEP SIX

K-Means Clustering
of Feature Set

STEP SEVEN

Review Cluster
Groupings

STEP EIGHT

PROVIDE RECOMMENDATIONS

PROSPECTIVE LOCATION CITIES

21 Moderate to large cities were identified within the South Eastern Region of the United States as prospective location cities. They are as follows:

- RALEIGH, NC
Pop: 479,332
- CHARLOTTE, NC
Pop: 859,035
- ATLANTA, GA
Pop: 486,290
- CHARLESTON, SC
Pop: 138,036
- WINSTON SALEM, NC
Pop: 244,605
- AUGUSTA, GA
Pop: 197,166
- COLUMBIA, SC
Pop: 133,114
- DURHAM, NC
Pop: 267,743
- RICHMOND, VA
Pop: 223,170
- NORFOLK, VA
Pop: 244,703
- ALEXANDRIA, VA
Pop: 160,035
- FLORENCE, SC
Pop: 38,317
- SAVANNAH, GA
Pop: 146,444
- ASHEVILLE, NC
Pop: 91,902
- ROANOKE, VA
Pop: 99,837
- KNOXVILLE, TN
Pop: 187,347
- GREENSBORO, NC
Pop: 290,222
- CARY, NC
Pop: 165,904
- ATHENS, GA
Pop: 125,691
- ALBANY, GA
Pop: 73,801
- CHATTANOOGA, TN
Pop: 177,571

DATA SOURCES

Information for Target and Prospective cities is compiled by webscraping source html code, and API queries with the python request library from the sources below.

- **WORLD POPULATION REVIEW**
Web scraped HTML resource providing Population Denisty information
www.worldpopulationreview.com
- **WIKIPEDIA**
Web scraped HTML resource providing Population information
www.wikipedia.com
- **AREA VIBES**
Web scraped HTML resource providing General Demographic data
www.areavibes.com
- **FOUR SQUARE**
Venue Search API based on Latitude/Longitude Location
www.foursquare.com

COMPILED DATA

Webscraping and API queries produces a list of compiled demographic and venue category information for each city. Compiled variable data points are displayed below:

WORLD POPULATION REVIEW

POPULATION DENSITY km/sq

WIKIPEDIA

LATITUDE

LONGITUDE

ESTIMATE POPULATION

2010 CENSUS POPULATION

AREA VIBES

CRIME PER 100 K

VIOLENT CRIME PER 100 K

PROPERTY CRIME PER 100K

INCOME PER CAPITA

FOUR SQUARE

VENUE NAME

VENUE LATITUDE

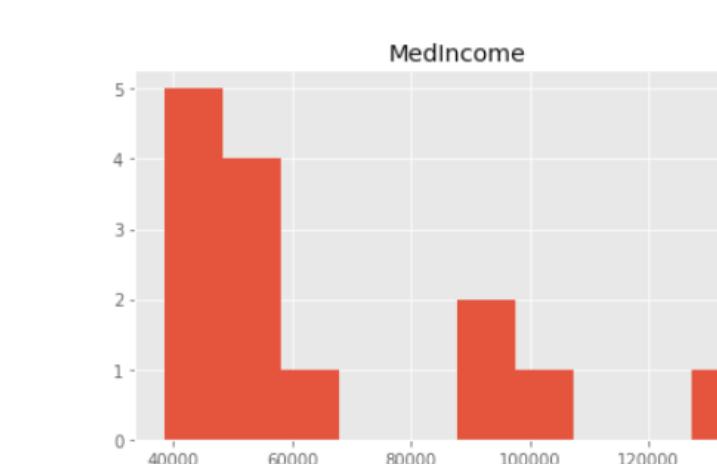
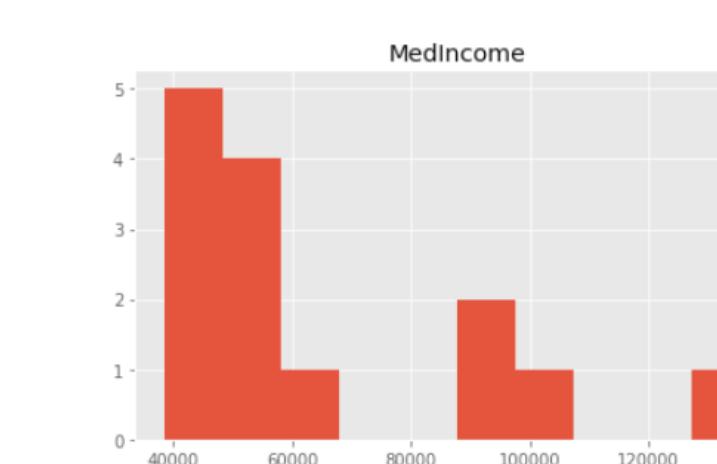
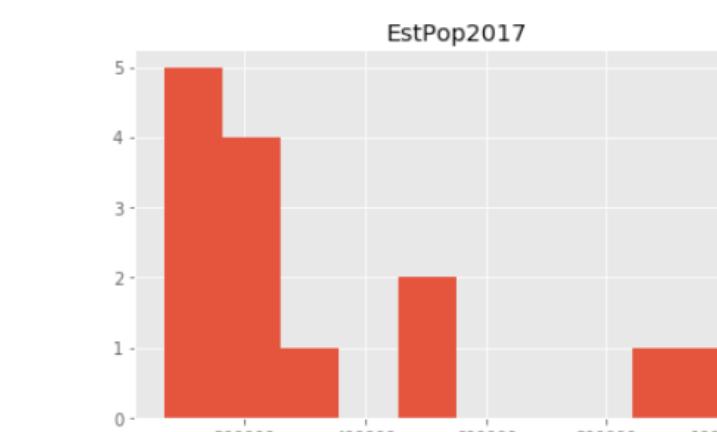
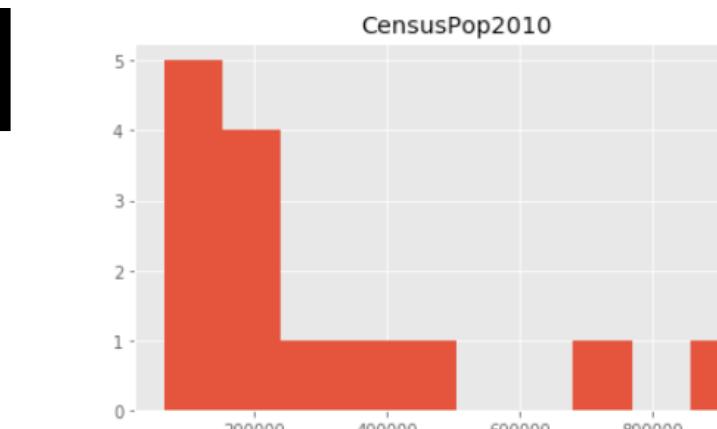
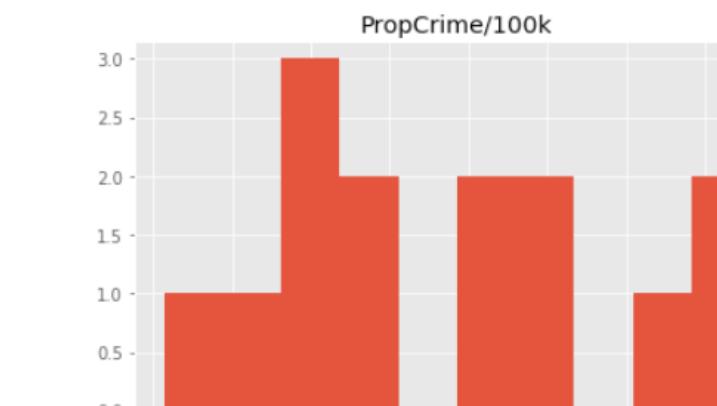
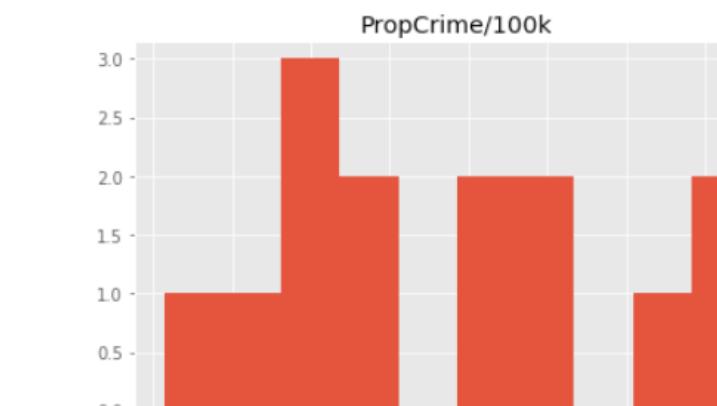
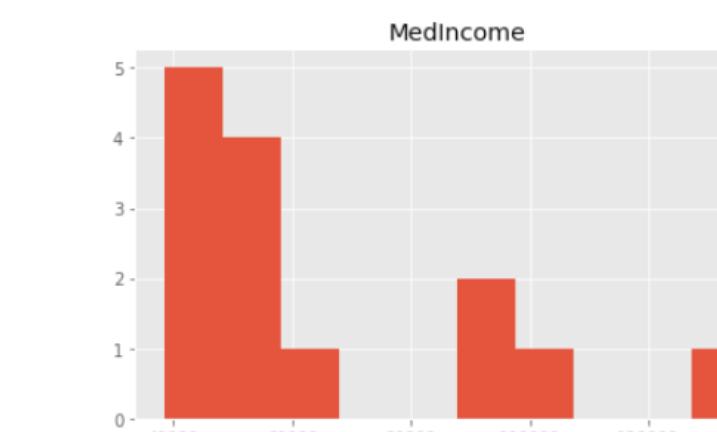
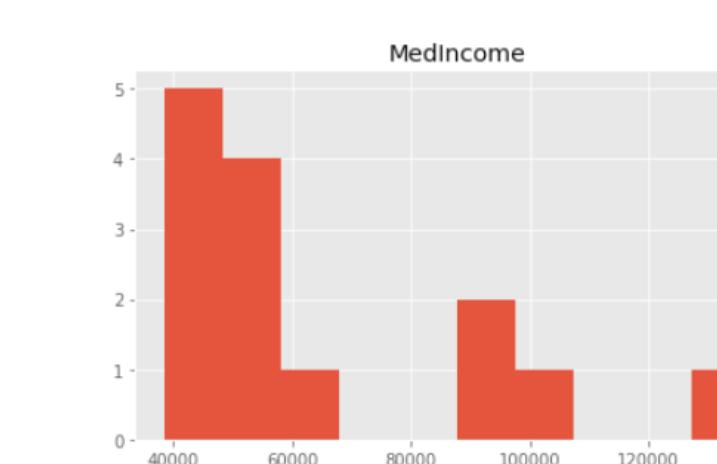
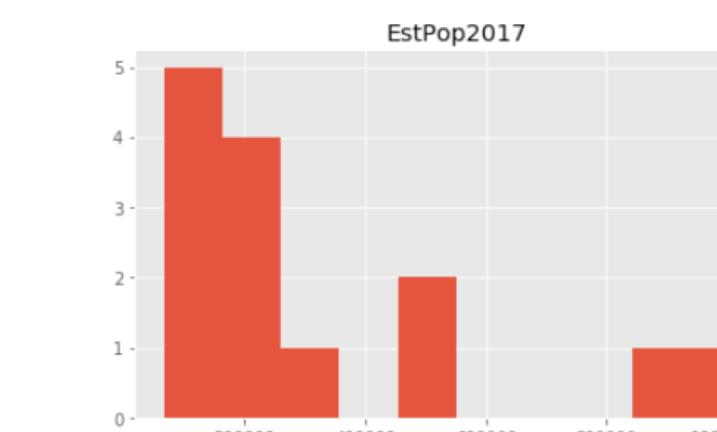
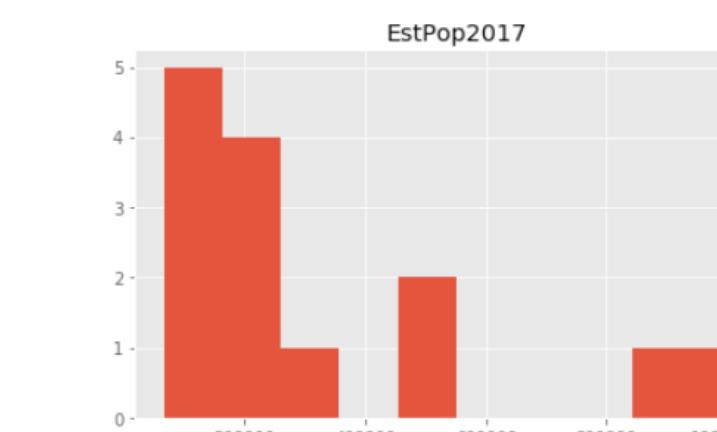
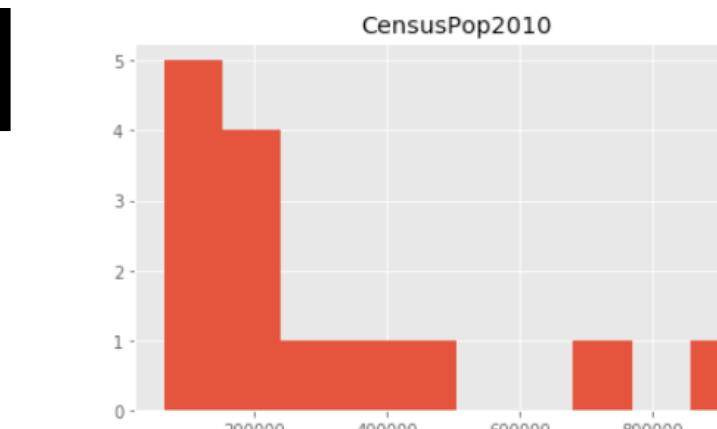
VENUE LONGITUDE

VENUE CATEGORY

DATA DISTRIBUTION ANALYSIS

Review compiled demographic variable data distribution for target, control and prospective cities to better understand the structure of compiled data.

The set of compiled variables show a diverse set of distributions with some uniform, right and left skewed distributions. Where Population and Med Income were the most skewed.



CORRELATION ANALYSIS

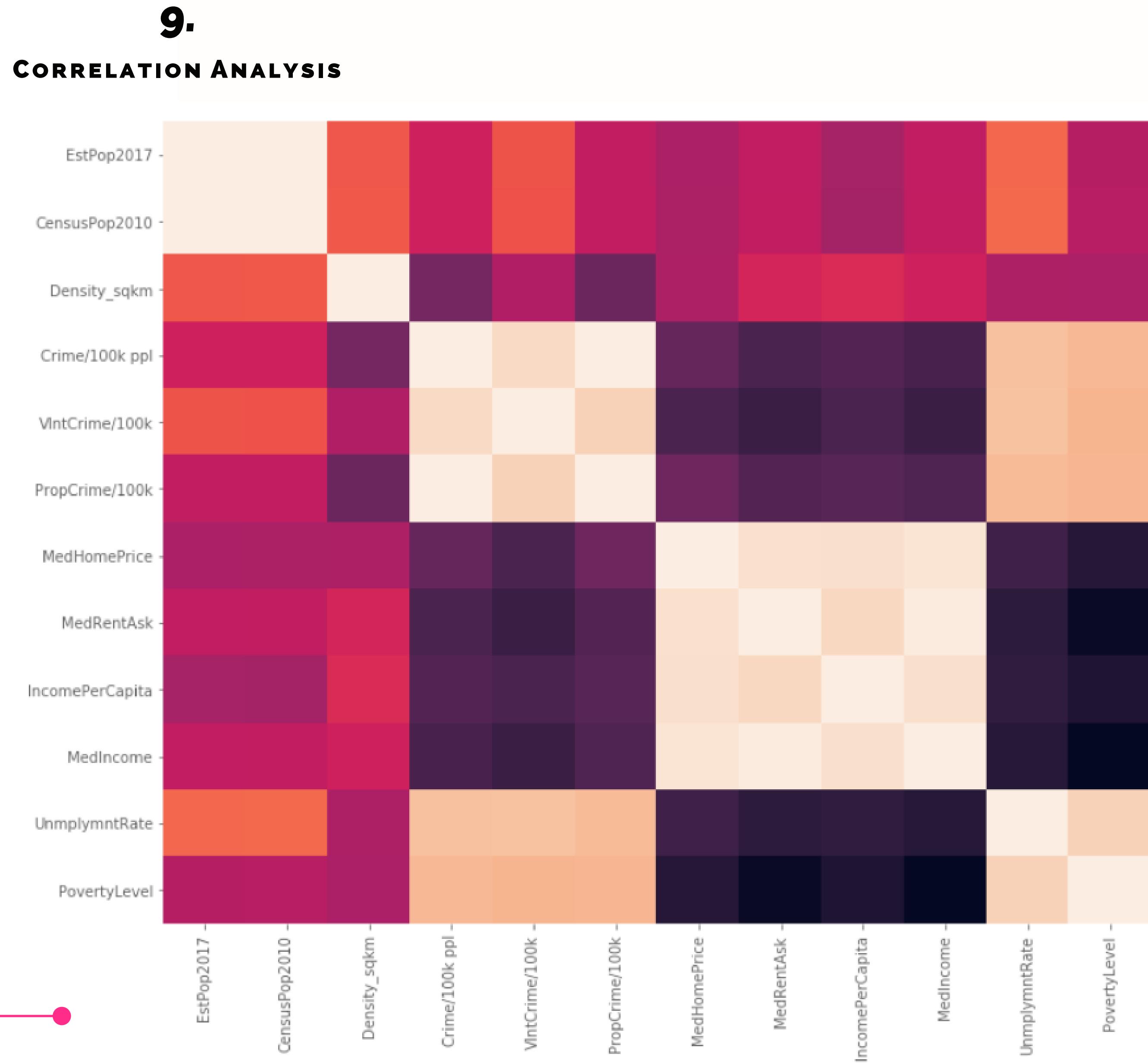
Review compiled demographic variable correlation to make feature determinations.

Strong Positive Correlations:

- Population, and 2010 Census Population
 - Crime per 100k people & Violent Crime per 100k/ Property Crime per 100k

Strong Negative Correlations:

- Poverty Level and Median Income
 - Poverty Level and Median Rent Ask



FEATURE SELECTION

Review of visualization and analysis led to Estimate Population and Census Population to be replaced by Growth Rate. Crime per 100k was removed from the data set. Selected demographic features are scaled with the min/max standard scaler and merged with venue category data.

GROWTH RATE

POPULATION
DENSITY km/sq

VIOLENT CRIME
PER 100 K

PROPERTY CRIME
PER 100K

INCOME PER
CAPITA

MEDIAN INCOME

UNEMPLOYMENT
RATE

POVERTY LEVEL

VENUE CATEGORY
AVERAGE

Growth Rate:

Derived by subtracting 2010 Census Population from Estimate Population, and dividing the value by Estimate Population

Crime p/100k: Removed

Venue Category Average:

Venues are Onehot encoded for the Category variable, The sum of onehot category count is divided by total venues. The count proportion is derived for each category

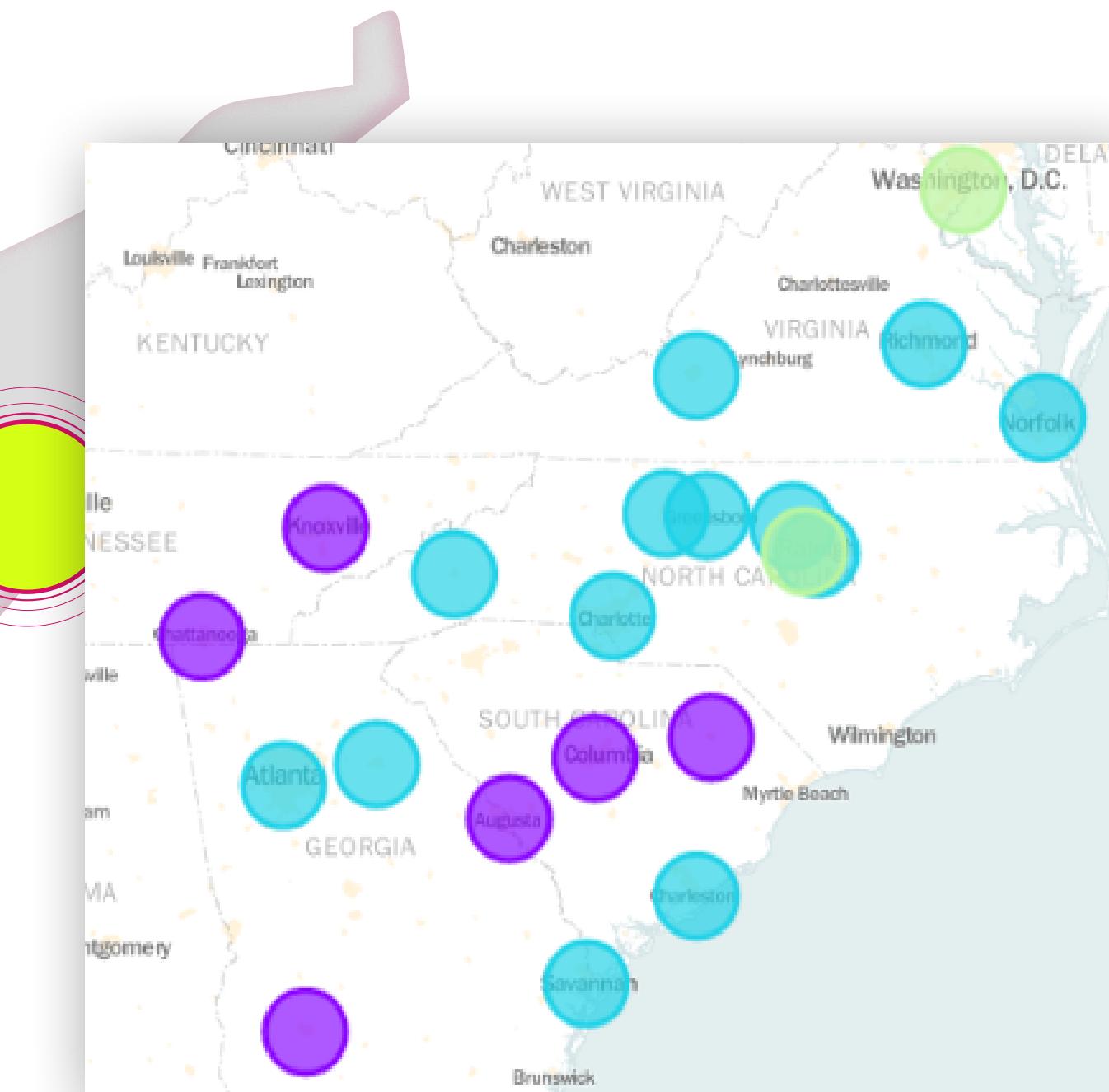
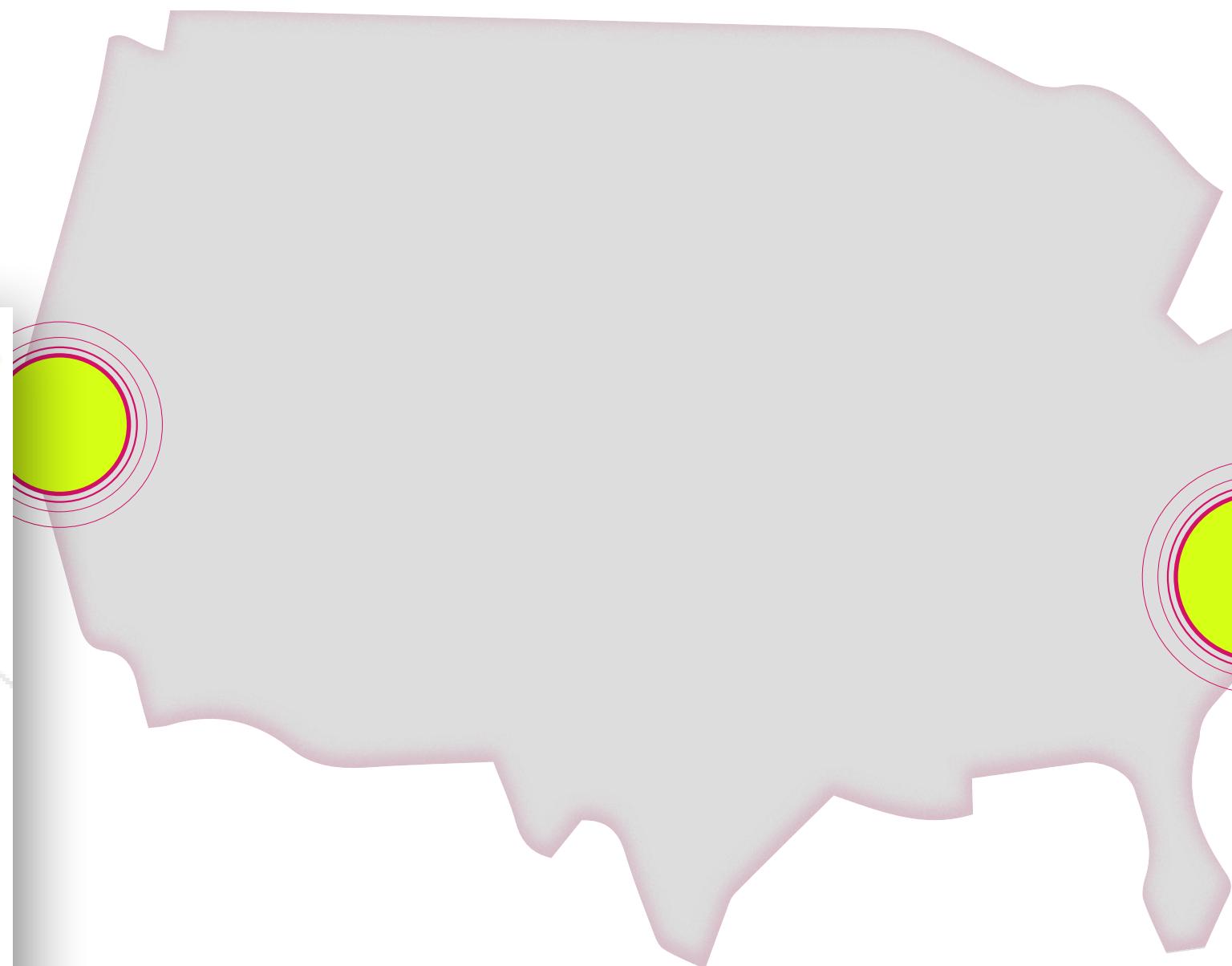
K-MEANS CLUSTERING

Given the small sample set cluster size is set to 3 distinct groupings with a random starting point value of 0. The K-Means Clustering algorithmn assigned groupings as follows:



Target Group Clustering:
Cluster Assignments are visualized with the Folium map plotting library.

Prospective Location Reccomendation:
Cary, NC or Alexandria, VA



RECOMMENDATIONS

With all features and venue categories considered the Cary, NC and Alexandria, VA cities would be the recommendation for prospective locations similar to Silicon Valley.

Note: Additional iterations of the K-means algorithm were conducted to test the recommended results with 2 & 4 possible groupings. In both instances Cary, NC and Alexandria, VA were clustered with the Target City.

CONSIDERATIONS

- Control Cities with varying demographics were also grouped with the Target city and Recommended Cities.
- Given the small size of the sample set K-means was set to group cities into 3 clusters as the ideal scenario
- There are a large number of additional demographic data points that could also be included that may affect clustering results.
- Feature selection and scale strategy may affect clustering results.

14.
CONTACT

TIMOTHY BRANNON

Raleigh, NC 27610 | brannon.timothy@gmail.com