

# **AFFORDABLE HOUSING DEVELOPMENT**

---

Group 4

Elio Aybar, Martin Copello,  
Matthew Fligiel, Matt Norgren

# Executive Summary



Affordable housing developments aid cities in promoting diversity and inclusion, improving quality of life, enriching neighborhoods, and growing the local economy. In this study, we investigate whether the locations of affordable housing developments provide positive features for its population, relating to accessibility to public transportation, amenities, active commercial scene, walkability, and general neighborhood activity.

Through our investigations, we discovered stark discrepancies between neighborhoods with a multitude of affordable housing developments and those without. Generally, areas with scarce access to transit and worse socio-economic performance contain most affordable housing developments.

# Project Objective

This project seeks to develop an end-to-end data pipeline for affordable housing data to:

- Examine data between the location, price, and size of affordable housing developments and standard-of-life features available for tenants.
- Foster a better understanding of:
  - Factors that affect location, prices, amenities, and quality of affordable housing developments.
  - The distribution of socio-economic indicators, economic activity, access to transit and parks, in each neighborhood.
- Potentially identify correlations that may confirm or deny the team's initial hypotheses related to affordable housing.



# Business Use Cases



Understanding affordable housing development's characteristics' relationship to socio-economic features could be used to develop more conscious policies and drive significant investments that provide inhabitants of affordable housing a better quality of life (i.e. development of additional public transportation in isolated neighborhoods, investment in local businesses, etc.)

# Data

## Data sources included:

- City of Chicago: Affordable Housing Development Dataset, Business Licenses Dataset, Park Dataset, Abandoned Property Datasets
- Zillow and Cook County Records API
- Neighborhood Shapefiles and Boundary files for GIS-type applications including generation of Table Keys (Neighborhoods), Public Transit and City Park GIS data points.
- Web-Scraping Walk-Scores of Chicago neighborhoods

## Data Profile:

- 428 data points, 75 columns, 85.5K rows





# Tools

Ingestion and  
Cleanup



Storage



Delivery and  
Insights



- Data was scraped via Python (Walk-Scores), downloaded from online sources (City of Chicago for Housing Developments, Commerce, Park) and gathered using APIs (Google Geocoding, Zillow, Cook County).
- Pandas was used for data cleanup.
- R was used to build an additional Geocoding API (Google's) to obtain coordinates from address.
- Our data database was built using MySQL. DDL and DML tables were created for each table.
- Tableau was used to generate tables and charts to draw insights from our data.

# Data Cleanup

## Out-of-scope Geographically

A	B	C	D	E	F	G	H	I	J
ID	LICENSE	ACCOU	SITE NU	LEGAL	DOING	ADDRES	CITY	STATE	ZIP CO
1888878-2	2570948	21545		1 PSF MECH	PSF MECH	11621 E	SEATTLE	WA	98168
1888878-2	2705036	21545		1 PSF MECH	PSF MECH	11621 E	SEATTLE	WA	98168

## Overwrote Assessed Value with Sale Value (where available)

```
72 df.loc[df['most_recent_sale_price'] == ('nan' or 'NaN')] = 0
73 df.most_recent_sale_price
74
75 #printing header names not values
76 def f(row):
77     if 'most_recent_sale_price' != ('nan' or 'NaN'):
78         'most_recent_sale_price'
79     else:
80         ('pri_est_land'+ 'pri_est_bldg')
81
82 df['price'] = df.apply(f,axis=1)
83
84 df.price
85 df['pri_est'] = (df.pri_est_land + df.pri_est_bldg)
86 df.pri_est
87 del df['pri_est_land']
88 del df['pri_est_bldg']
```

# GIS Conversions

```
from geopy.geocoders import Nominatim
locator = Nominatim(user_agent="myGeocoder")
ln = locator.geocode("455 Cityfront Plaza Dr, Chicago, IL 60611")
print(ln.latitude, ln.longitude)
```

41.8903132 -87.62168978061763

```
In [10]: from shapely.geometry import Point
def ca_assign4(pt):
    for i, v in enumerate(CA_boundaries):
        if pt[0]>v[2]:
            continue
        if v[1]>pt[1]:
            continue
        if pt[1]>v[3]:
            continue
        if v[0]>pt[0]:
            continue
        elif CA_shapes[i].contains(Point(pt[0], pt[1])):
            return i
    break
```

```
In [11]: from shapely.geometry import Point
import pickle
CA_names = pickle.load(open("C:/Users/Martin/Desktop/CA_names.pkl", "rb"))
CA_boundaries = pickle.load(open("C:/Users/Martin/Desktop/CA_boundaries.pkl", "rb"))
CA_shapes = pickle.load(open("C:/Users/Martin/Desktop/CA_shapes.pkl", "rb"))
nbc_tow = (ln.longitude, ln.latitude)
CA_names[ca_assign4(nbc_tow)]
```

Out[11]: 'NEAR NORTH SIDE'

```
1 getwd()
2
3 # Geocoding a csv column of "addresses" in R
4
5 #load ggmap
6 library(ggmap)
7
8 # Select the file from the file chooser
9 fileToLoad <- file.choose(new = TRUE)
10
11 # Read in the CSV data and store it in a variable
12 origAddress <- read.csv(fileToLoad, stringsAsFactors = FALSE)
13
14 # Initialize the data frame
15 geocoded <- data.frame(stringsAsFactors = FALSE)
16
17 # Loop through the addresses to get the latitude and longitude of each address and add it to the
18 # origAddress data frame in new columns lat and lon
19 for(i in 1:nrow(origAddress))
20 {
21     # Print("working...")
22     result <- geocode(origAddress$addresses[i], output = "latlon", source = register_google(key = "AizaSyDe9Yqysf7_0pquD21
23     origAddress$lon[i] <- as.numeric(result[1])
24     origAddress$lat[i] <- as.numeric(result[2])
25     origAddress$geoAddress[i] <- as.character(result[3])
26 }
27 # Write a csv file containing origAddress to the working directory
28 write.csv(origAddress, "geocoded.csv", row.names=FALSE)
```

# Design Considerations I: Data

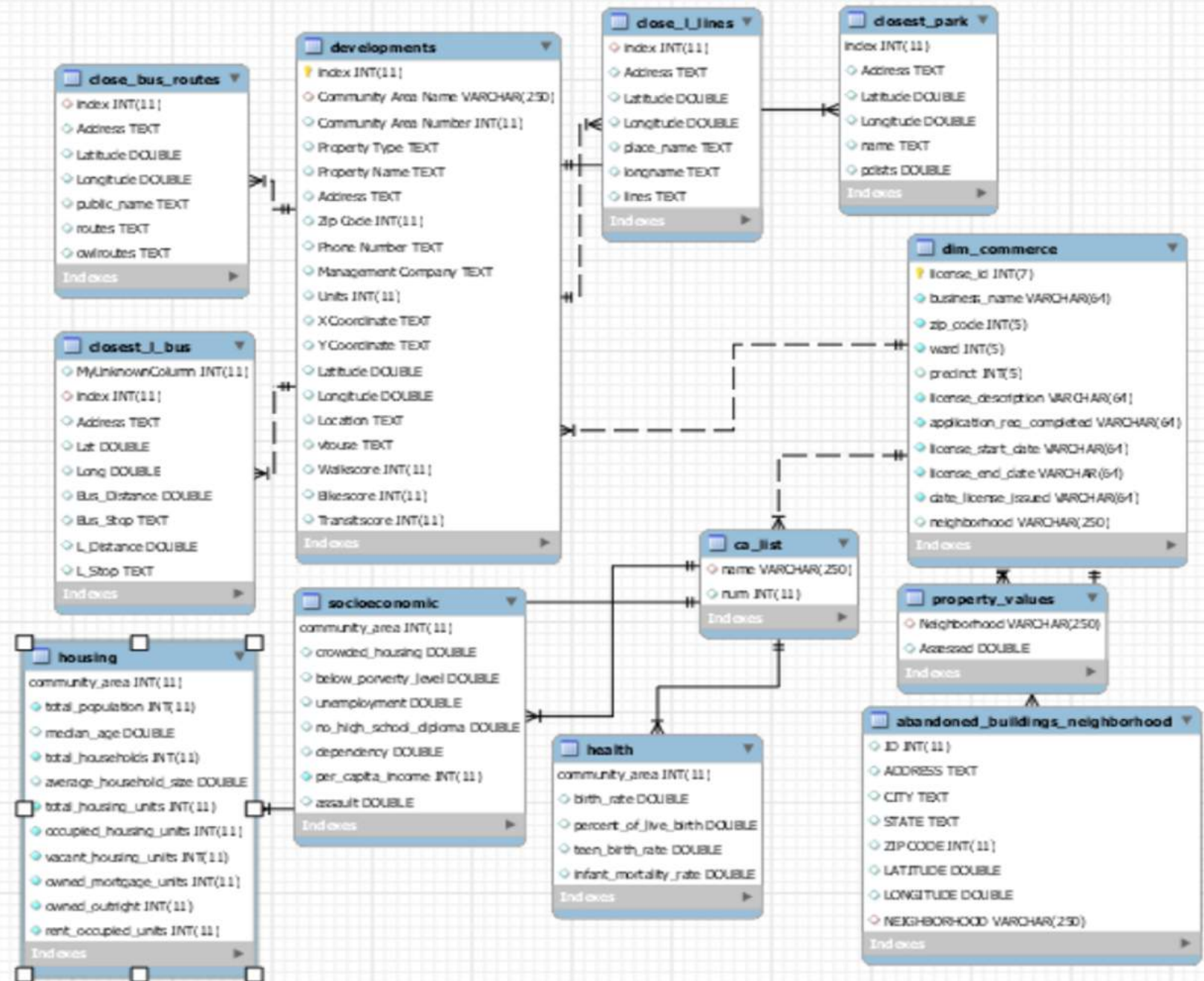
- ETL Scripts written in Python and R
- Keys were generated for tables (neighborhood from latitude and longitude, or from address):
  - Chicago neighborhood shapefiles
  - Geocodes via GeoPy and Google's Geocoding API via GCP to obtain Latitude and Longitude from address
  - Web-scraping Walk-Scores to calculate "walkability" of area



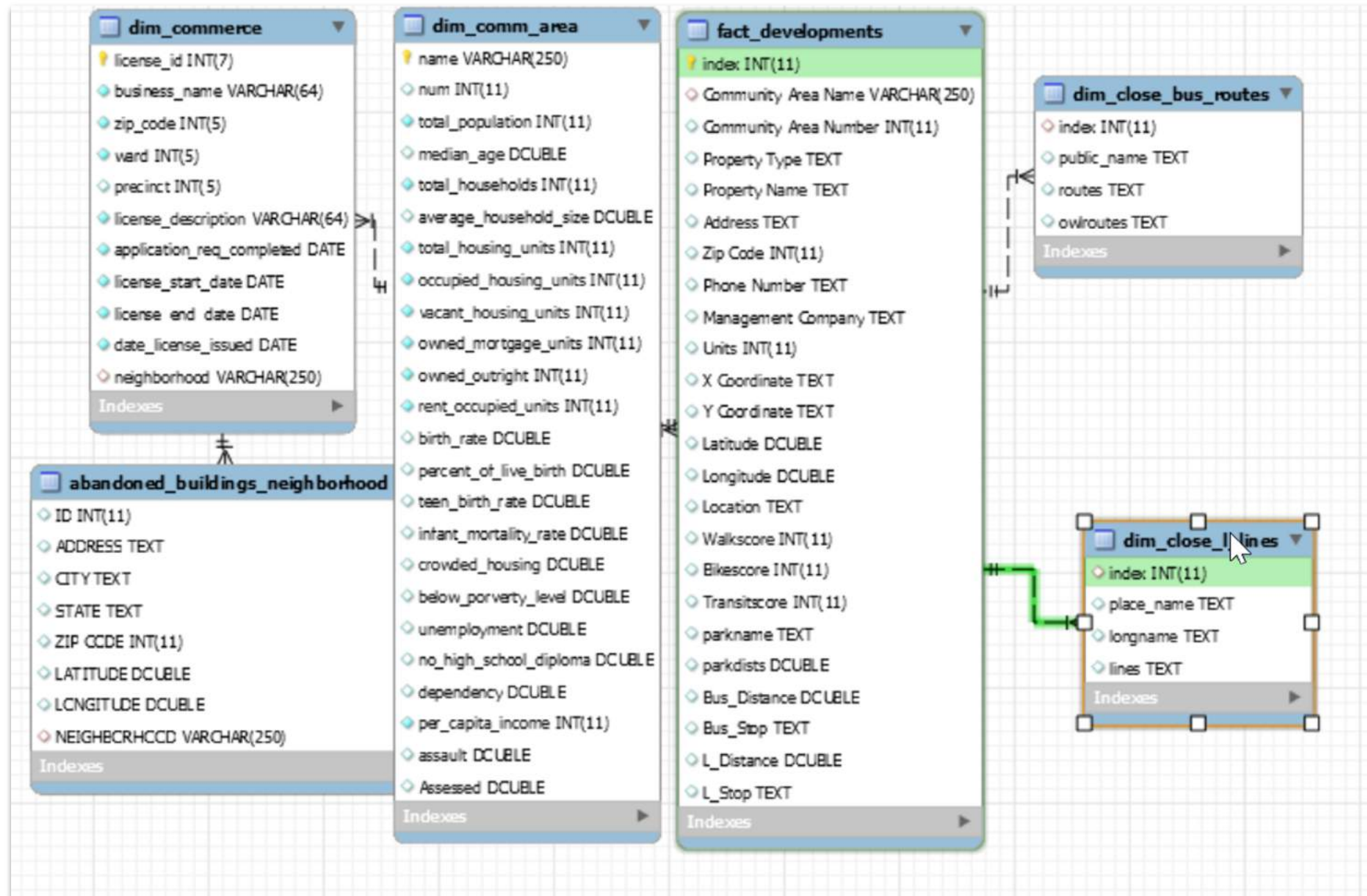
# Design Considerations II: Implementation

- Number of tables and joins are not extensive, so this pipeline can be implemented locally given enough disk space
- Due to the archival nature of our data, our application can be used as an OLAP system.
- Snowflake data model
- One to many relationship between individual address and neighborhood requires definition of granularity at neighborhood level

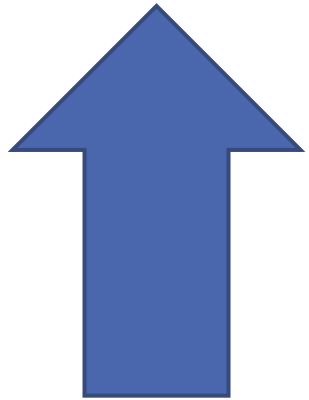
# EER Diagram



# Dimensional Model

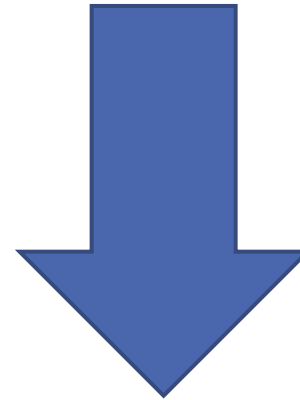


# Preliminary Expectations



**Greater number  
of affordable  
housing  
developments**

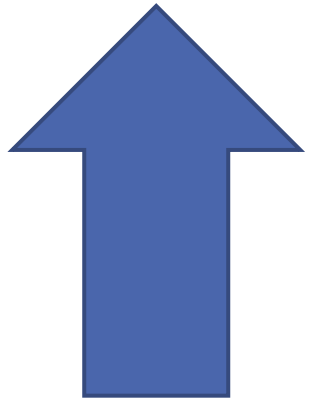
- In less affluent areas
- Near abandoned buildings
- In highly-commercialized areas



**Fewer number  
of affordable  
housing  
developments**

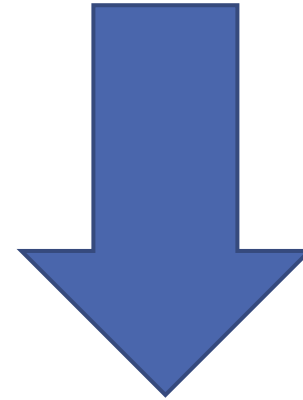
- Near public transit
- Near parks

# Actual Results



**Greater number  
of affordable  
housing  
developments**

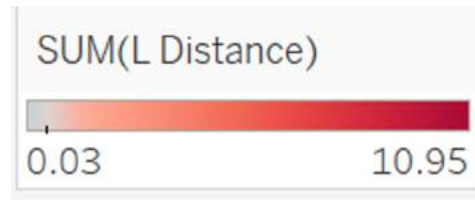
- In less affluent areas ✓
- Near abandoned buildings ~
- In highly-commercialized ~  
areas



**Fewer number  
of affordable  
housing  
developments**

- Near public transit ~
- Near parks X

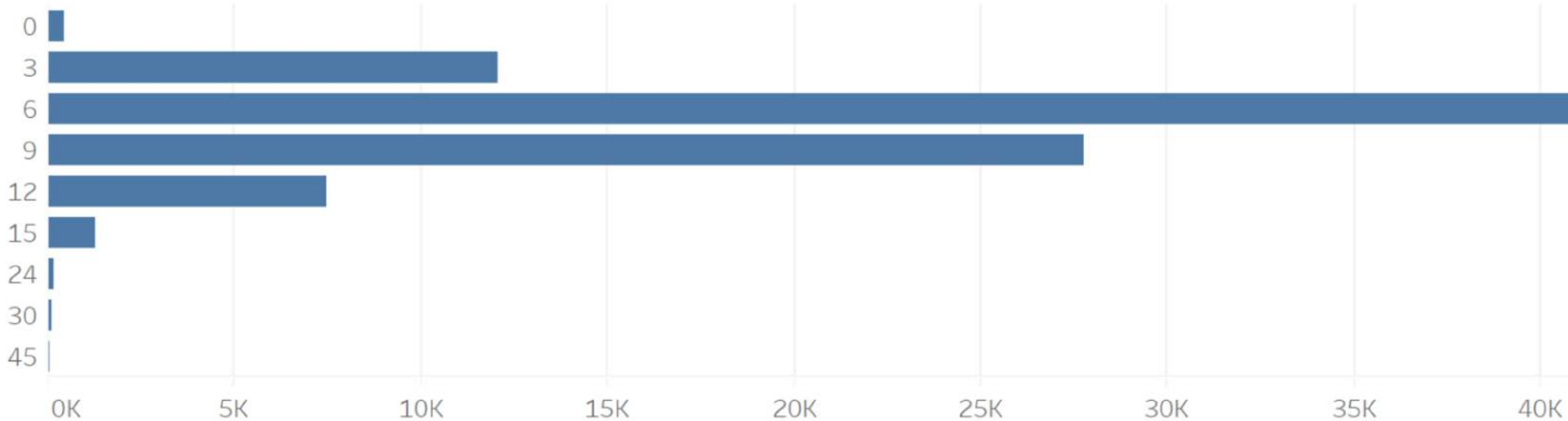




Distance from  
Housing  
Developments  
to One  
Instance of  
Public  
Transportation

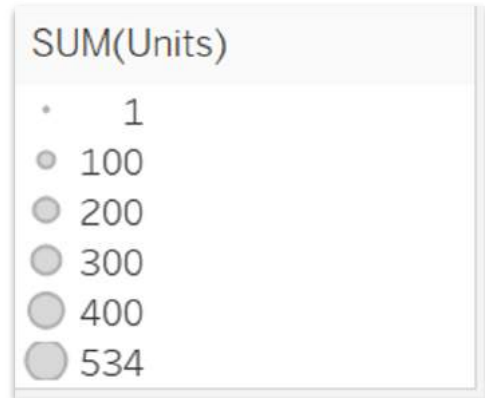
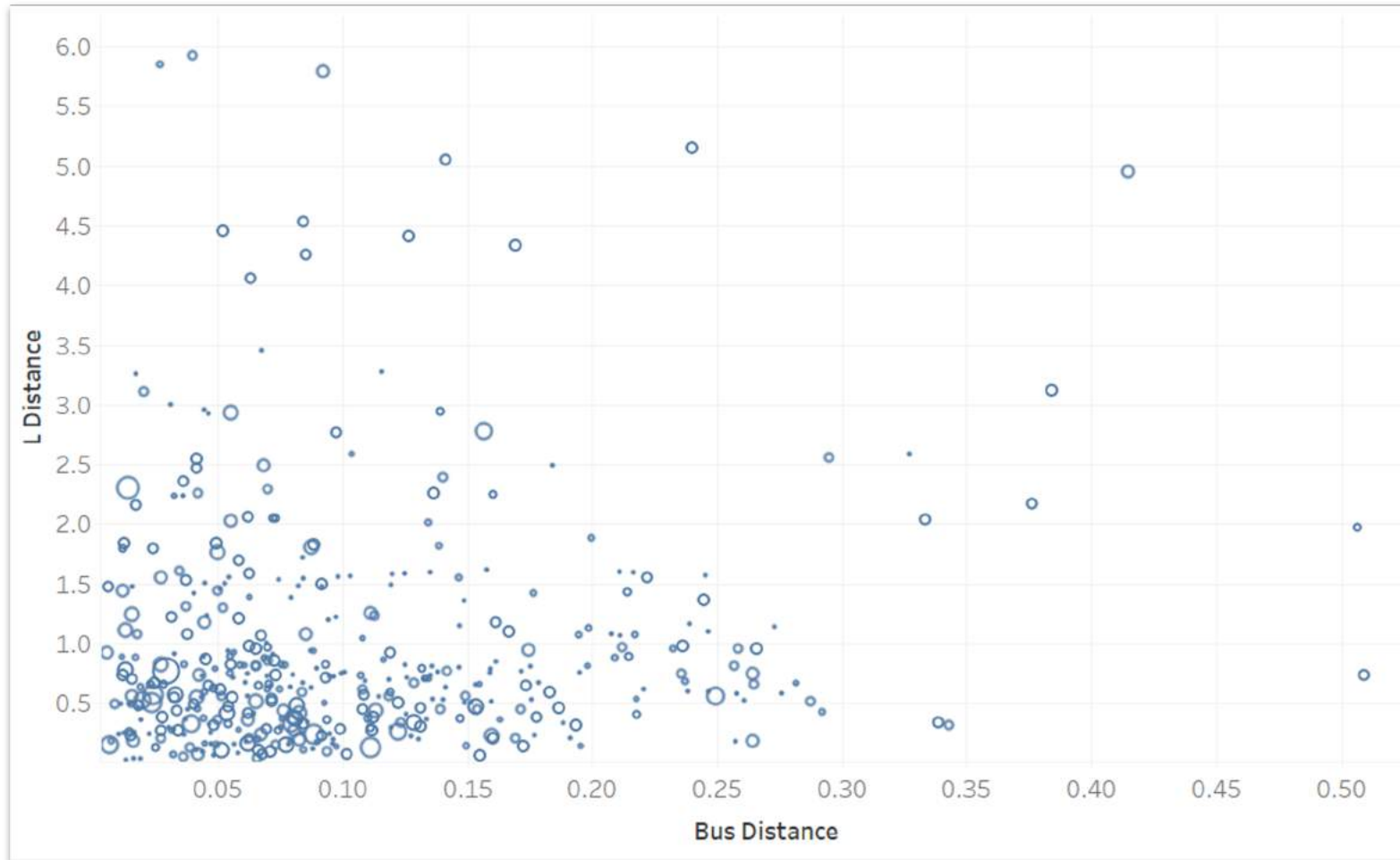
# Number of Bus Routes within 1-Mile of an Affordable Housing Development

CD\_Routes ..

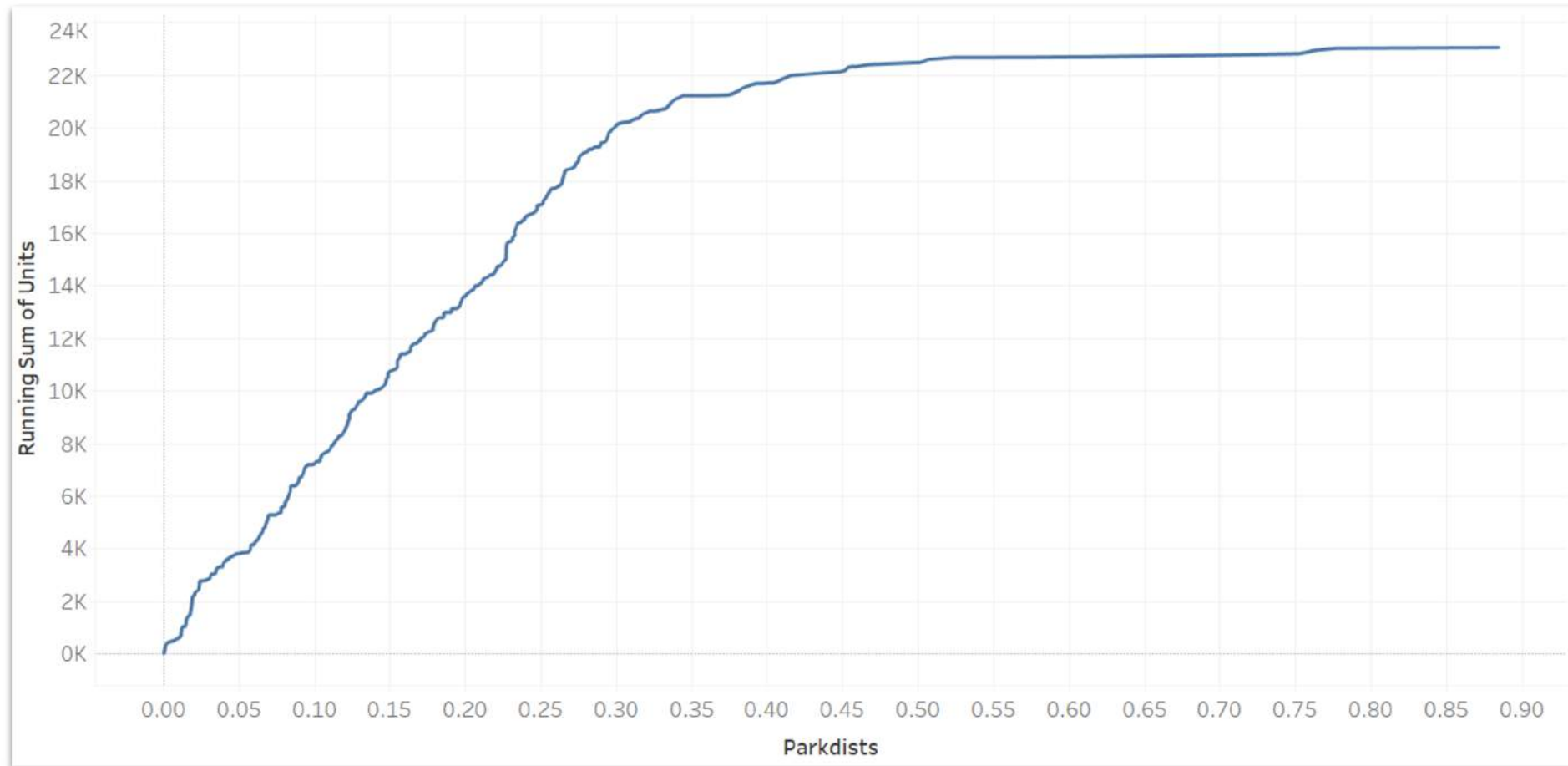


Index (Dim Close Bus Routes)

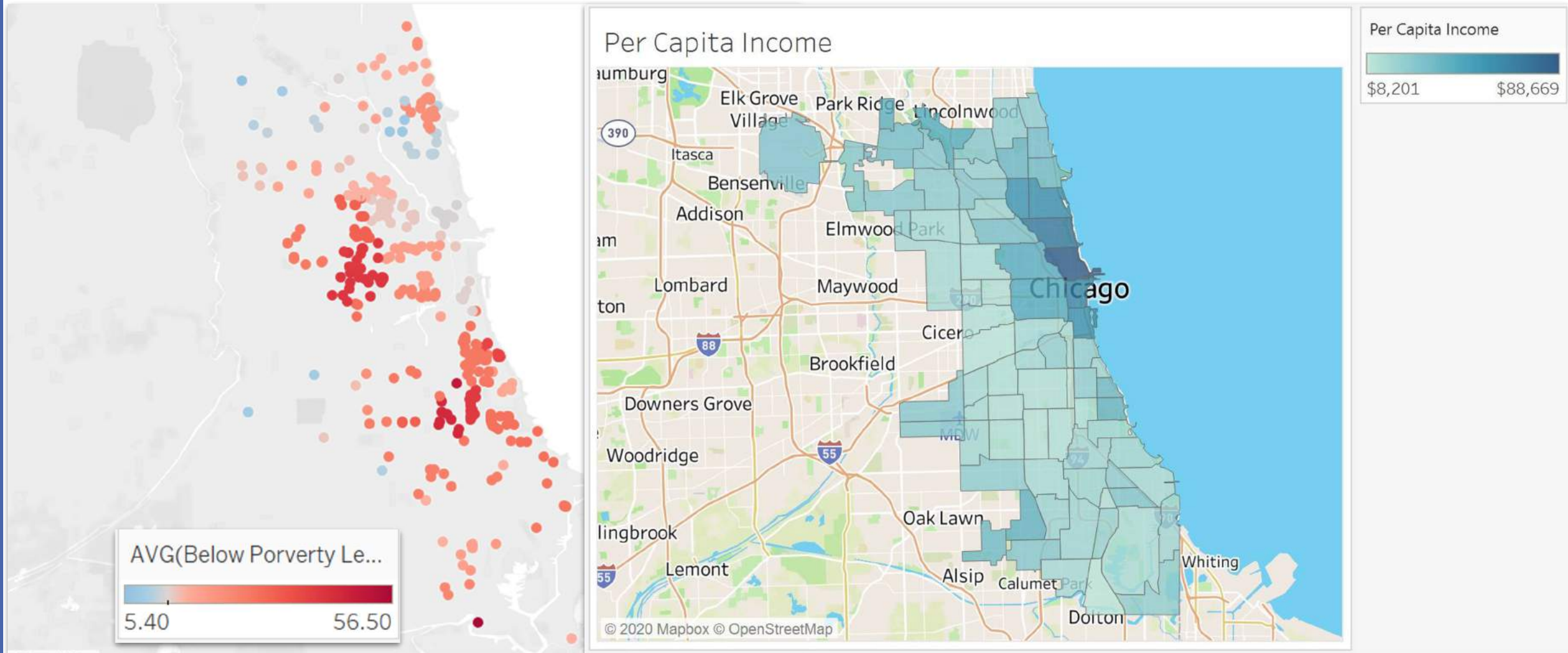
# Number of Units by Bus and L Distance



# Most units did have access to nearby parks

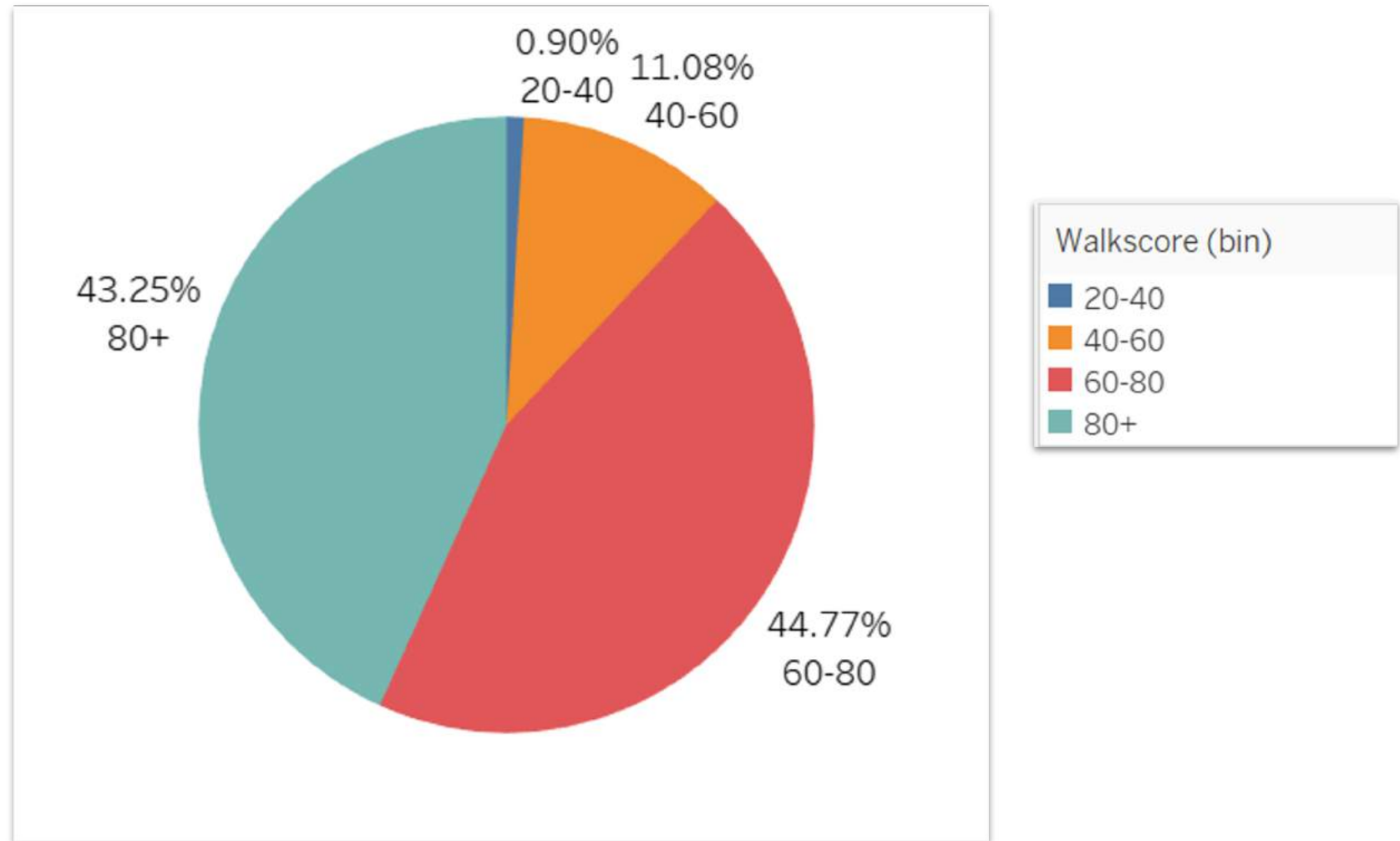


# Affordable Housing Developments and Percentage of Neighborhood Below Poverty Line (Centered at 13%)

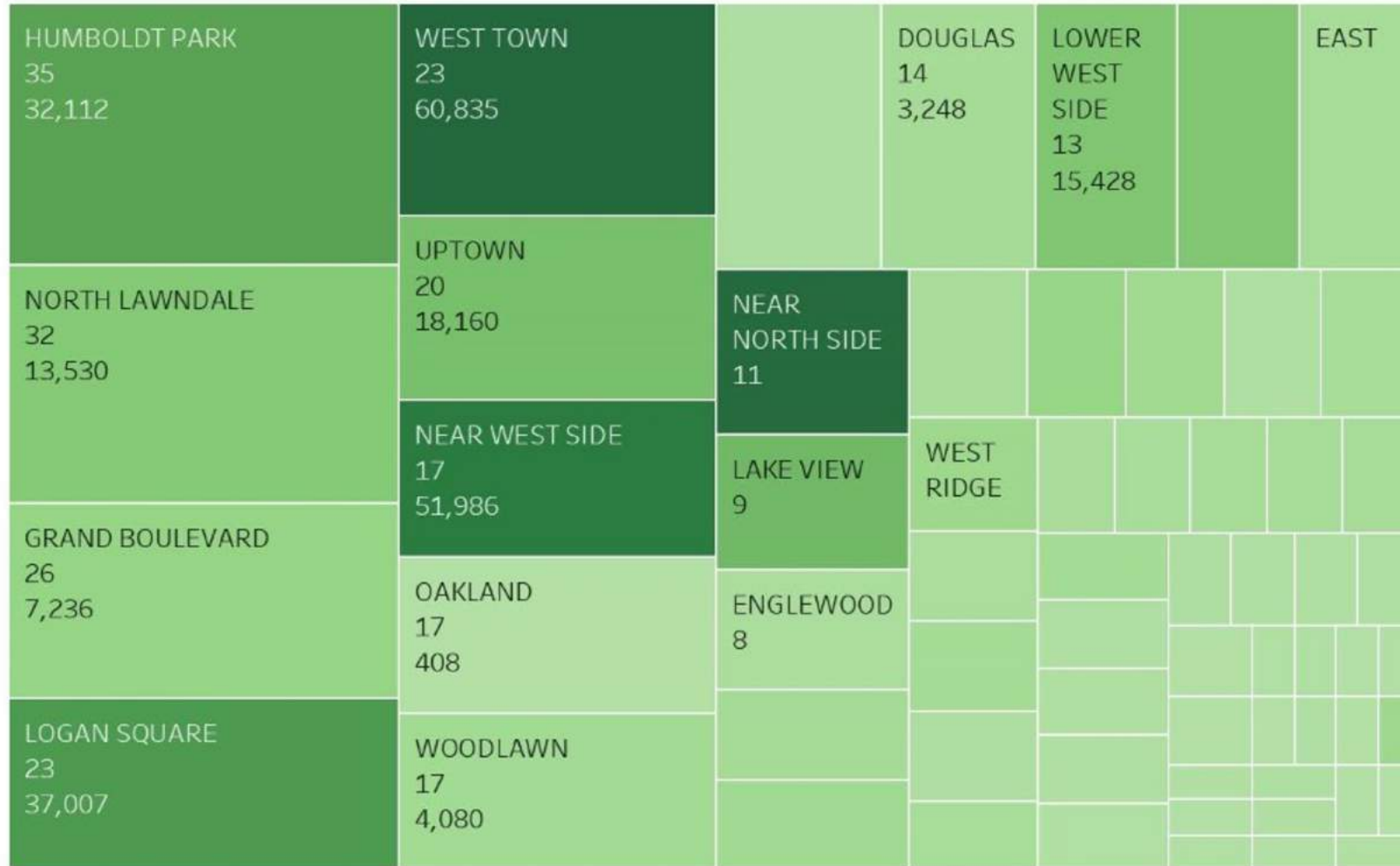




# Walkscore (Binned by # of Units)



## Housing Developments subset by Business Density

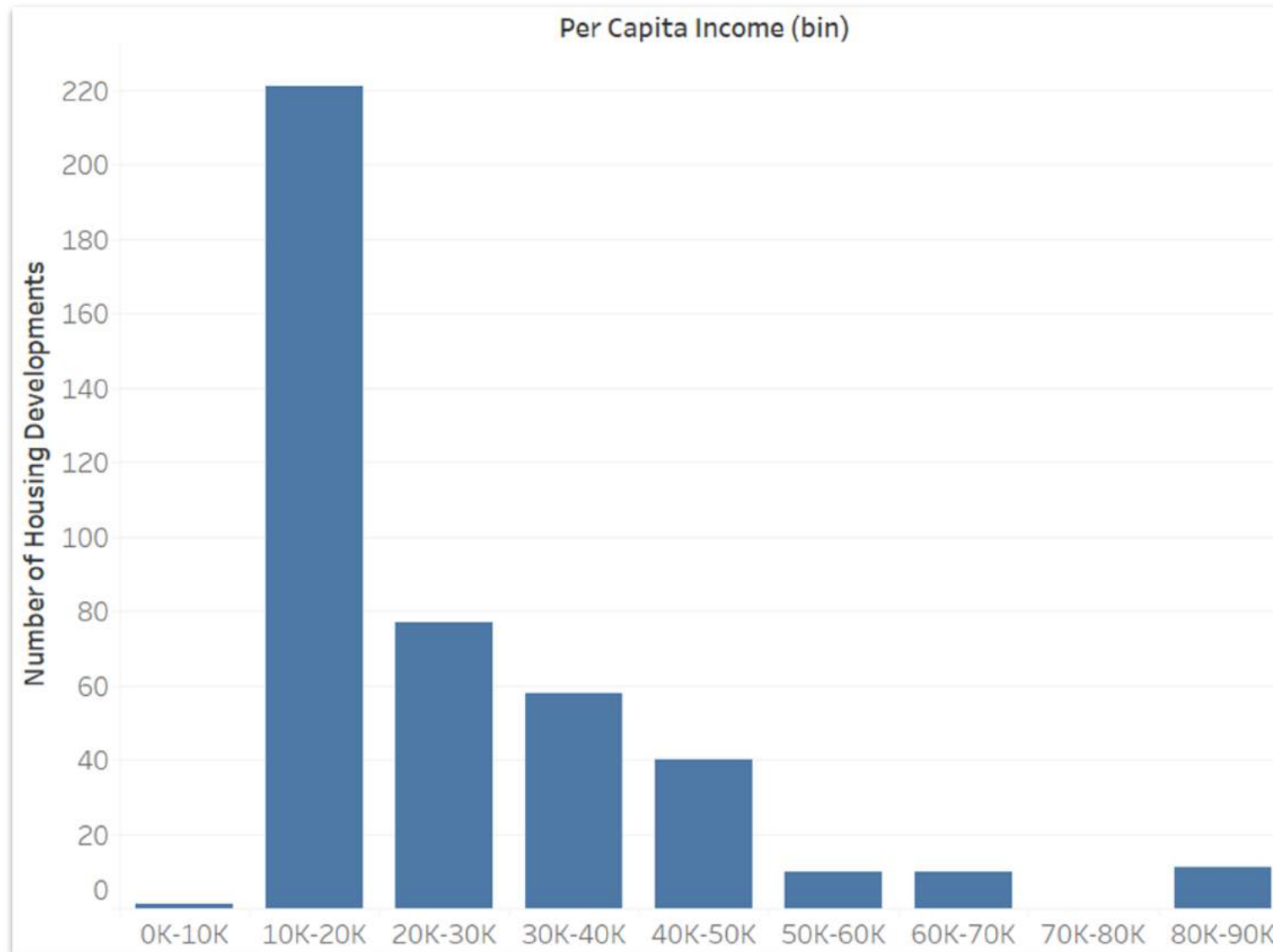


SUM(Businesses)

39

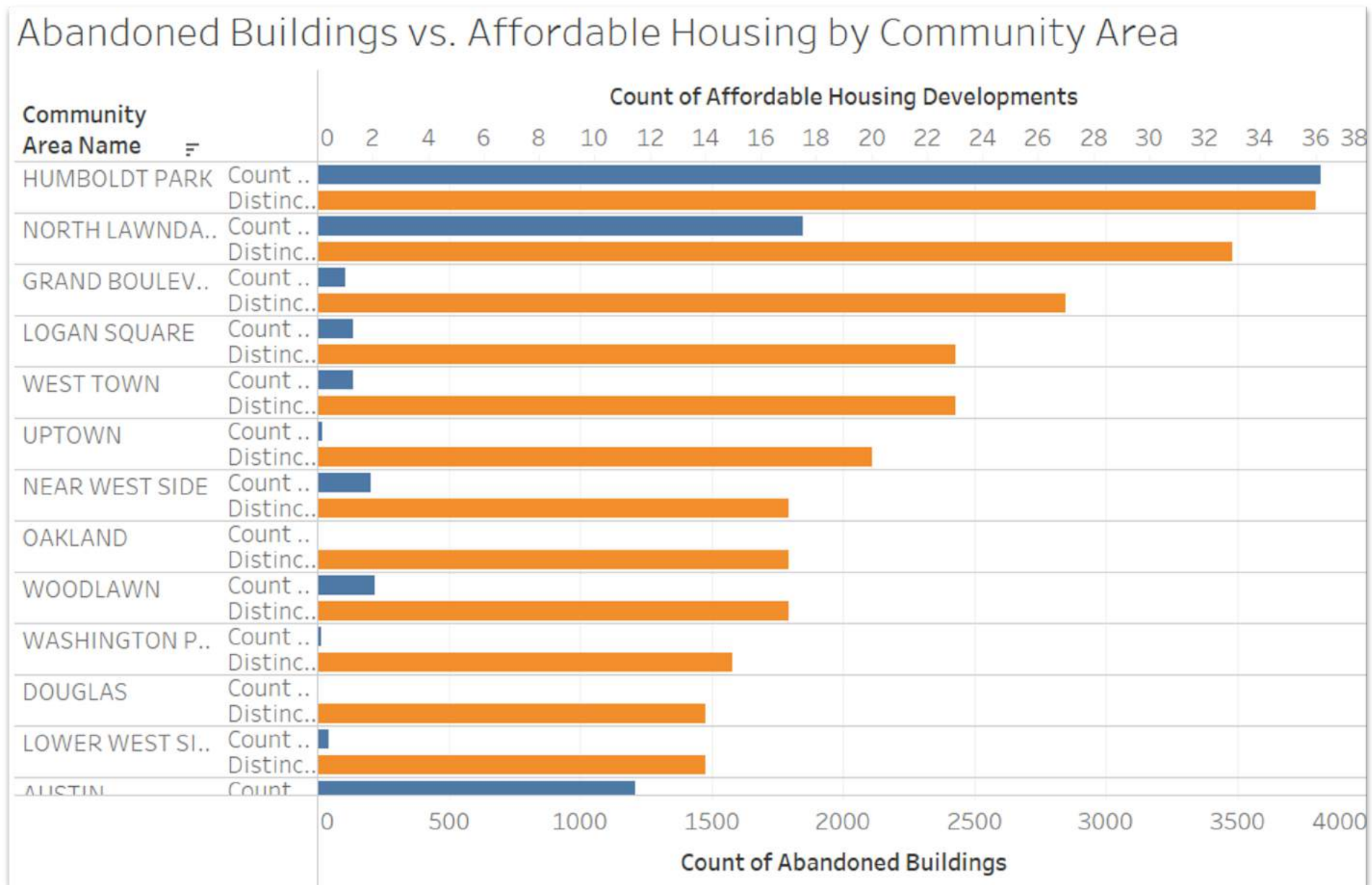
60,835

# Developments by Per Capita Income

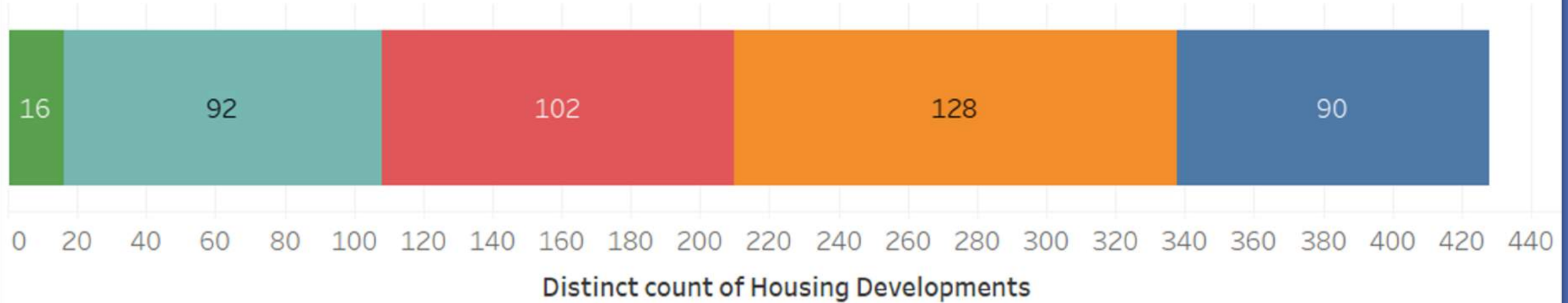


**Count of  
Affordable  
Housing  
Developmen  
ts**

**Count of  
Abandoned  
Buildings**



## Number of Developments by Household Size



Average Household

1.5-2

2-2.5

2.5-3

3-3.5

3.5+



# Conclusion

Through the course of the project, we discovered strong correlations between the count of affordable housing developments per neighborhood and:

- Socioeconomic status of neighborhood, which could perhaps be explained by higher land costs for more affluent neighborhoods
- Transit scores of neighborhood for L trains, which could perhaps further isolate neighborhoods with lower household income per capita.

Additionally, weaker correlations were identified between the count of affordable housing developments per neighborhood and:

- Business density of each neighborhood
- Count of abandoned buildings of each neighborhood

It is difficult to draw conclusions from these weaker correlations.

# Areas of Improvement for Project

- With enough computing resources and time, further granularity could be utilized to map by-block impacts, driven by the by address level development, assessment, and transportation data, rather than by the neighborhood or area level. This could yield further insight and inform policy by visualizing ripple effects related to affordable housing developments and the blocks in which they are built.
- Obtaining data from the 2020 U.S. Census and re-evaluate our findings
- Study additional datasets (type of housing, family size, health background, etc.)

# References

- City of Chicago Data Portal: <https://data.cityofchicago.org/>
- Institute for Housing Studies at DePaul University:  
<https://www.housingstudies.org/>
- Google Maps Platform:  
<https://developers.google.com/maps/documentation>
- Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython 1st Edition by Wes McKinney

