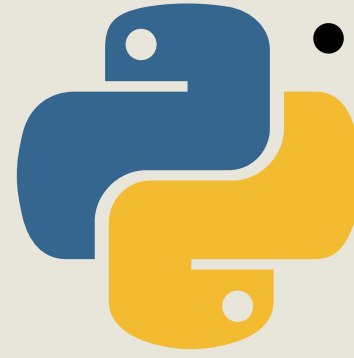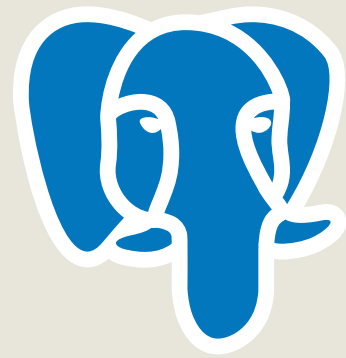# REAL ESTATE

## RADAR

By: Diahann Castellon: Data Engineer, Brandon Ingalz: Data Manager,
Anushya Mani: Data Architect, Roland "Coy" Abellano: Senior Data
Analyst, & Jimmy Nguyen: Data Scientist

# What City fits your Housing Budget & Needs?

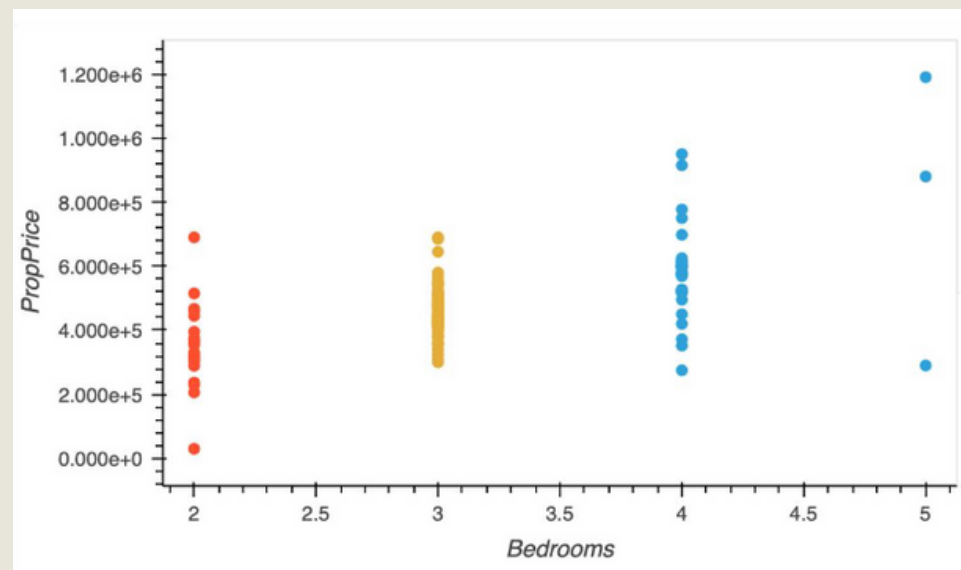| | Bathrooms | Bedrooms | Living Area | Land Area | Property Price | Inflation Rate | Federal Interest Rate | Month of the year | PCA Variance ratio (3 Components) |
|---|---|---|---|---|---|---|---|---|---|
| Chandler, AZ | X | X | X | X | X | | | | 0.967 |
| Chicago, IL | X | X | X | X | X | | | | 0.9979 |
| Los Angeles, CA | X | X | X | X | X | X | X | | 0.893 |
| Miami, FL | X | X | X | X | X | X | X | | 0.961 |
| New York, NY | X | X | X | X | X | X | X | | 0.983 |
| Portland, OR | X | X | X | X | X | | | X | 0.923 |

**Factors not included:**

- Real estate demand, per city/micro-city:
  1. Seasonal impact?
  2. Unemployment rates?
  3. Education rating?
  4. Safety?
  5. Political impacts?
- HOA fees per property
- Last renovation date, per property
- Walkability, per property
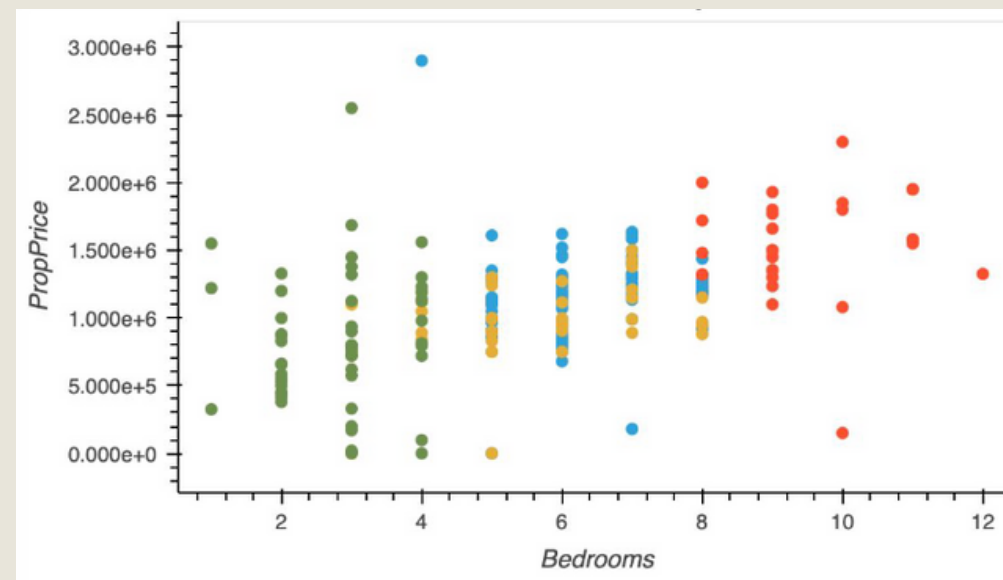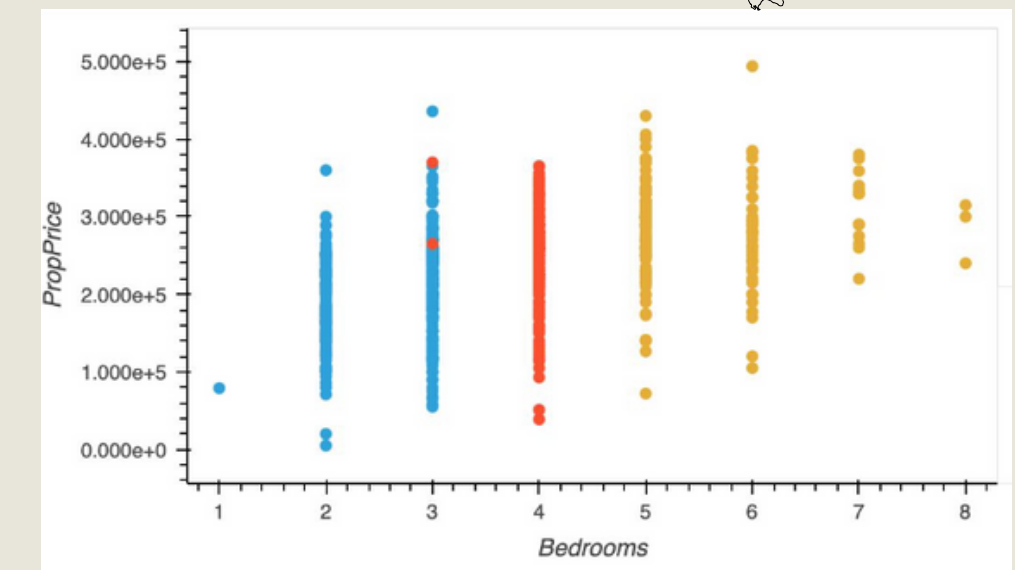- Other desirable property features such as: kitchen, basement, storage space, and parking

REAL ESTATE
RADAR

# Potential Application
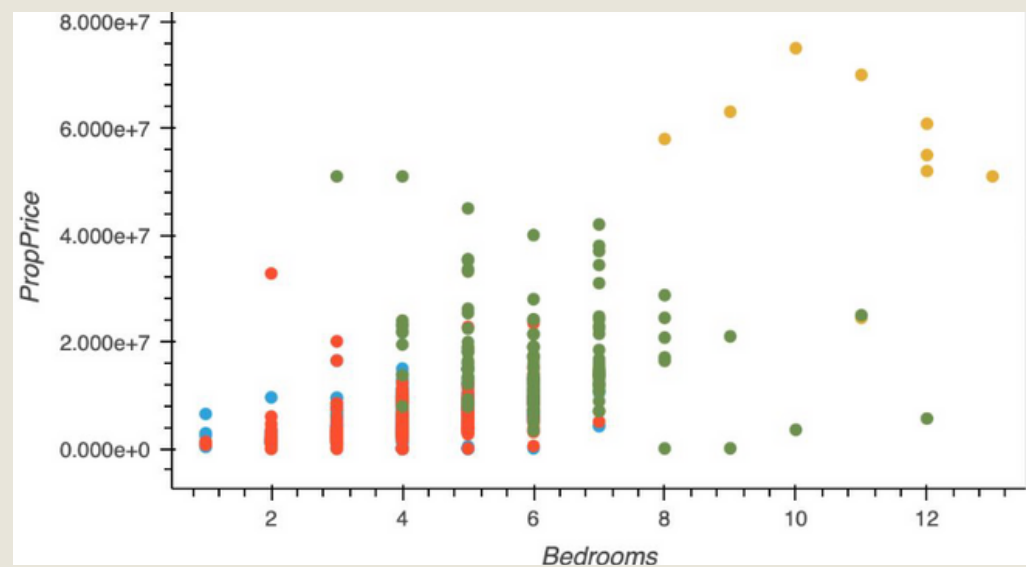
## Where would you buy next?

**Step 1:** Pick the housing cluster(s) you can relate to the most, per city

**Step 2:** Normalize the (max$/feature), per feature across clusters. Use the inverse of these values

**Note:** Remove any feature that did not prove to be significant in the cluster analysis!

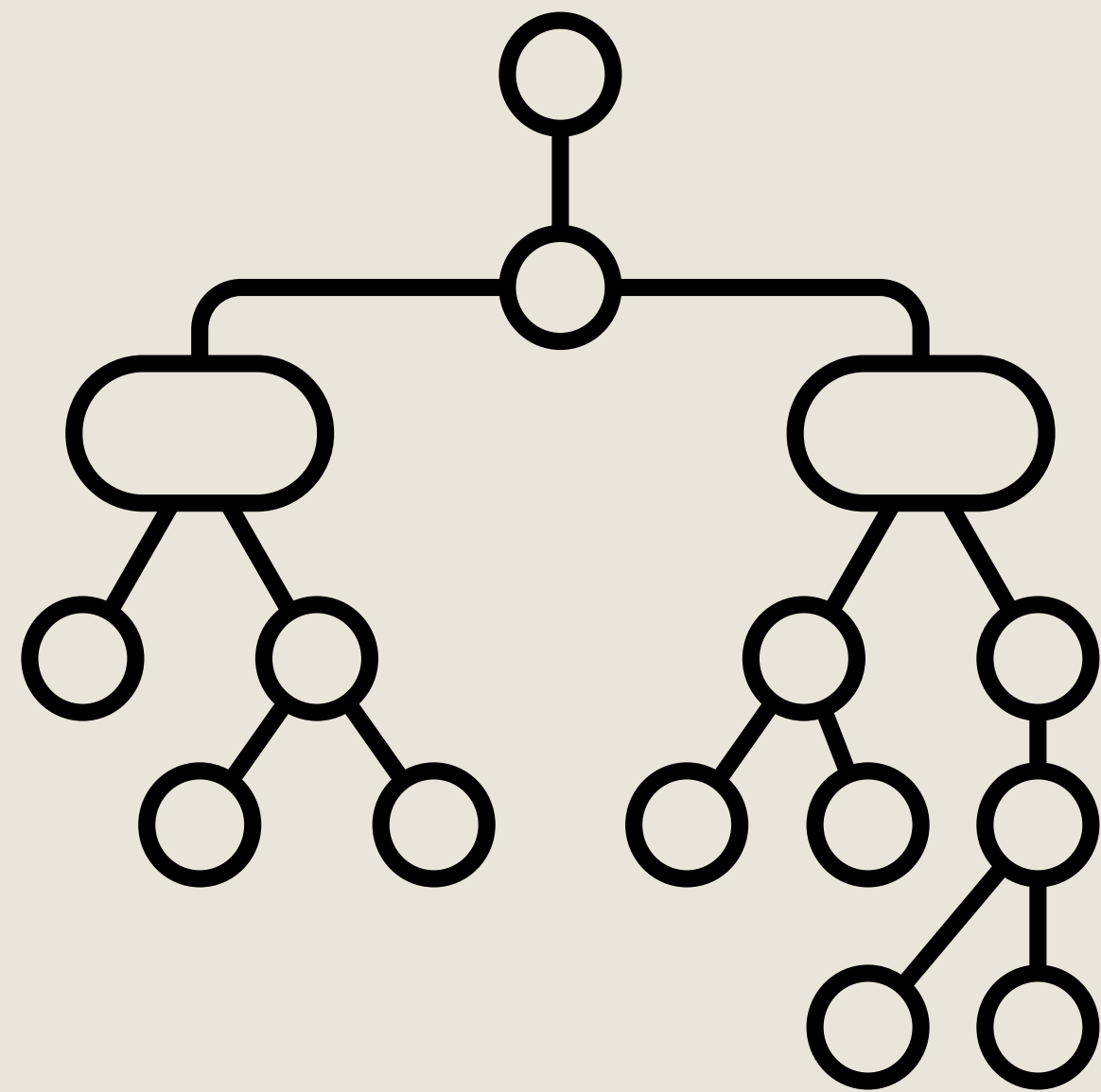**Step 3:** Determine how important each feature is to you

(Give a weightage for each property feature, sum should = 1):

**Step 4:** Calculate overall score for each cluster based on preferences, choose the cluster (and associated city) that scored the highest

(Example of scoring equation per city below):

**ClusterX_City Y** = (living_area_weight * ClusterX_CityY_norm["living_area] + Land_area_weight * ClusterX_CityY_norm ["land_area"] + ...)
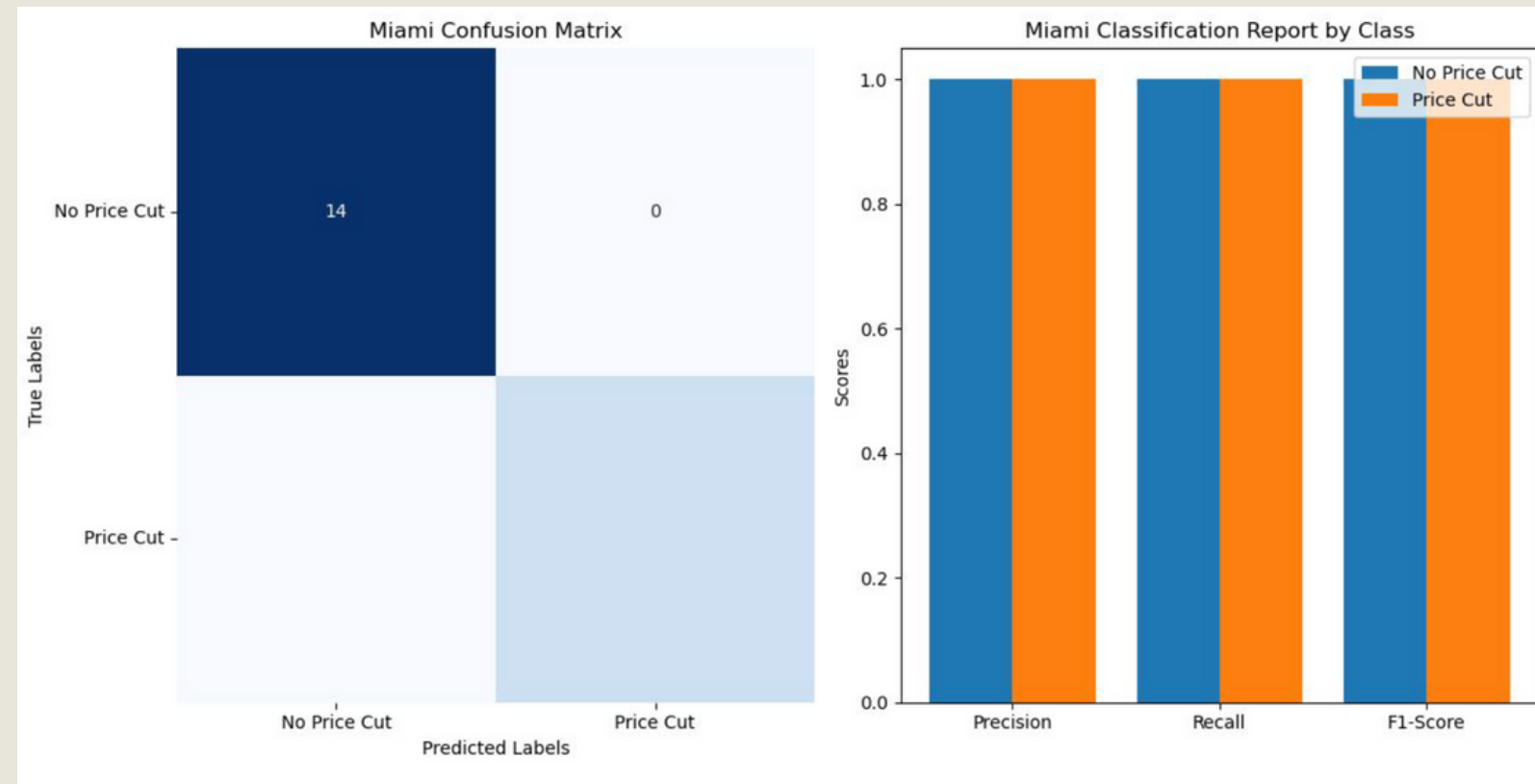
# Supervised ML Model
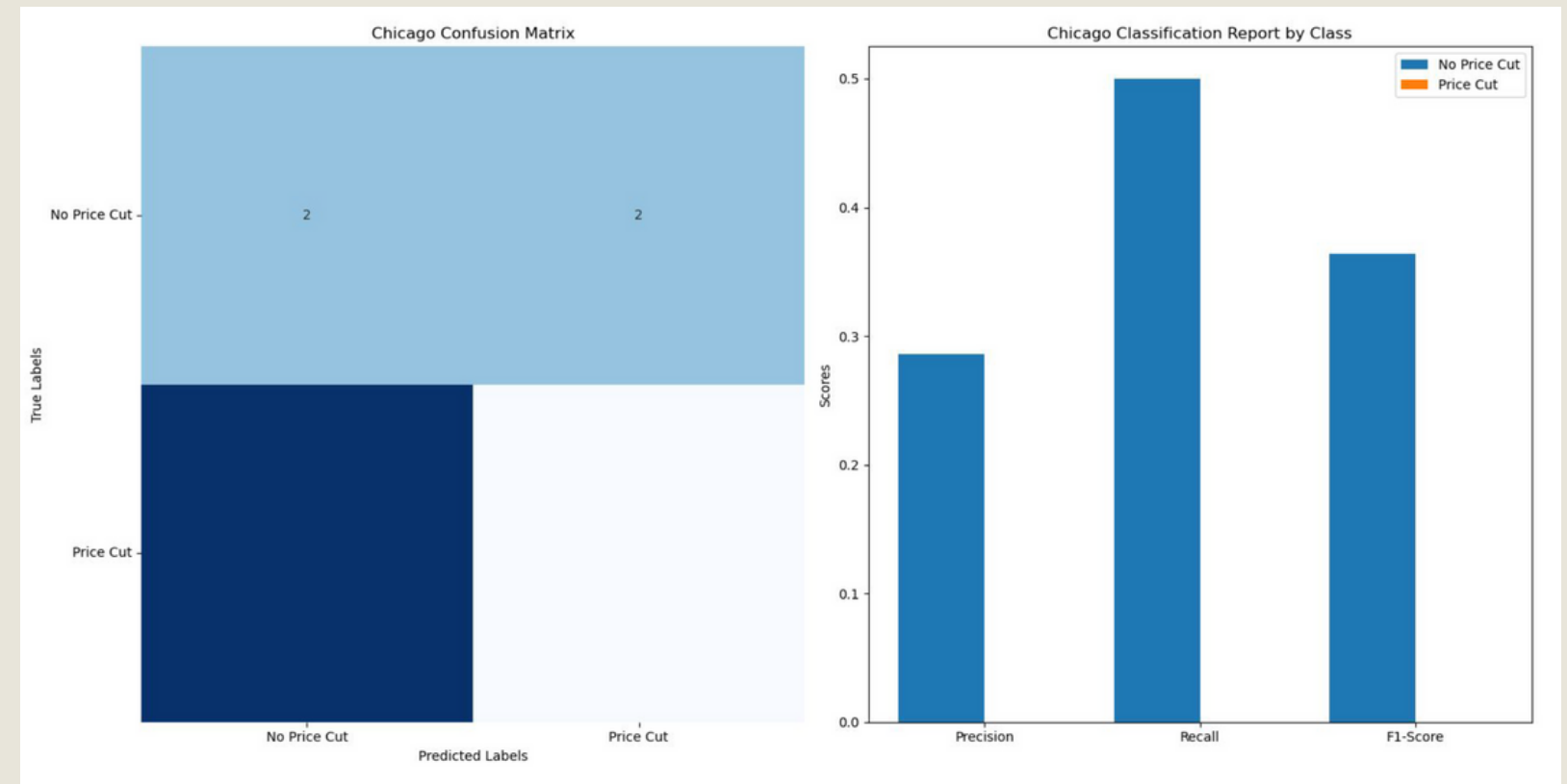
**New York City, NY**

**Los Angeles, CA**

# Supervised ML Model Continued
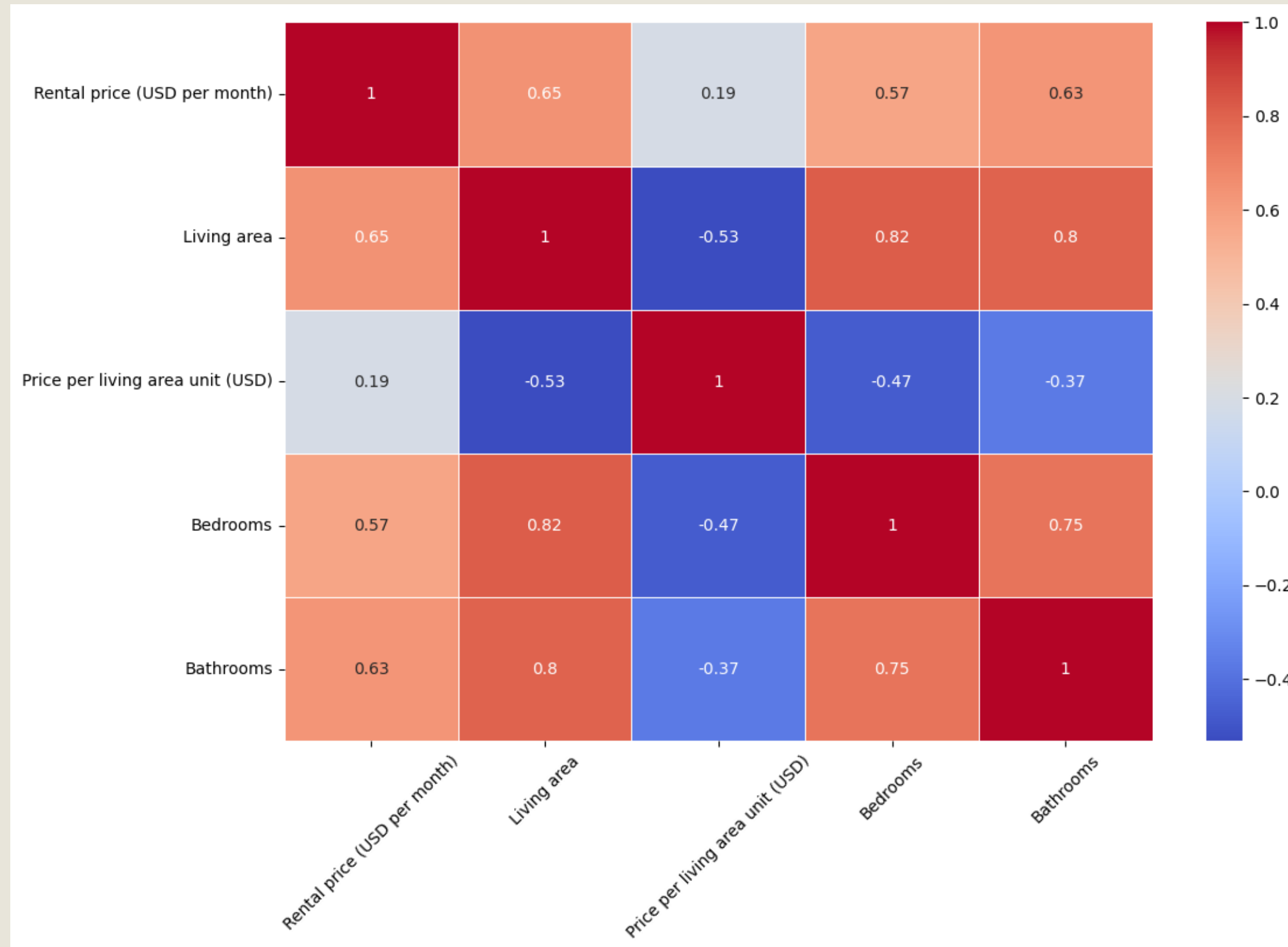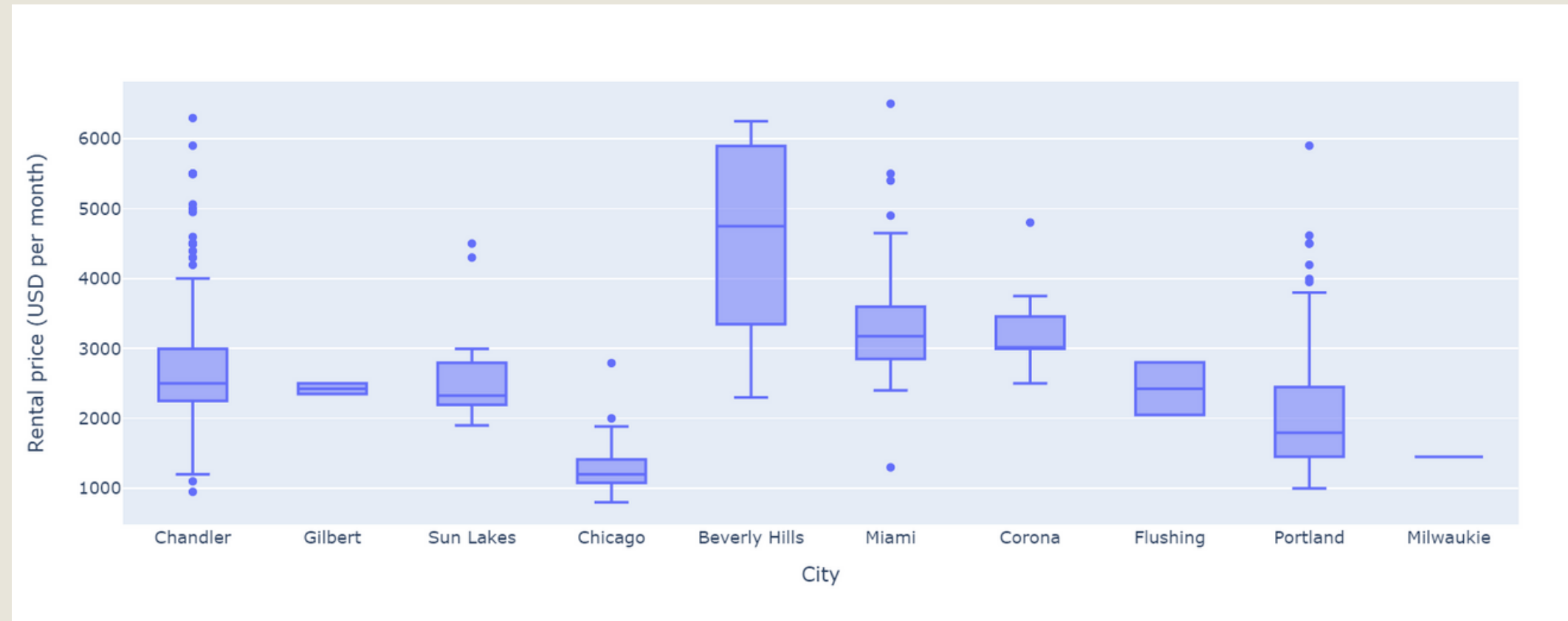
**Miami, FL**



**Chicago, IL**

# Correlation of Variables

Price vs Area

# Limitations

- Parts of the United States do not disclose property sales to the public
- Need to integrate additional features such as school ratings or transportation scores
- Data set limited to **36 months**