# ZILLOW DATA ANALYSIS: RECENTLY SOLD HOMES IN LOS ANGELES, CA (LAST 90 DAYS)

**Mohammad Amin Yousefi - Ali Hamzehpour - Mina Shirazi**

**-Spring 2024-**

# AGENDA

- **Introducing our dataset:**

    - **What's Zillow?**
    - **Why Zillow?**
    - **Why Los Angeles?**

- **Our Goal**

- **Scraping data**

- **Preprocess Steps**

- **Analyzing Distribution Outliers**

- **TOP 10s**

- **Findings from Visualizing Feature Changes Based on Year Built**

- **Correlation**

- **Group Bar Plots**

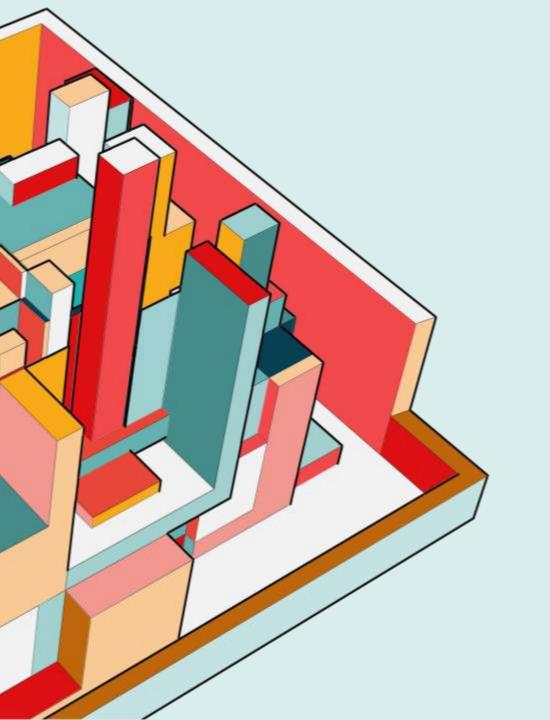- **Growth Rate**

- **Statistical Test**

- **Next Step…**

# WHAT'S ZILLOW?

Zillow is a popular online real estate marketplace that provides information about homes, real estate listings, mortgages, and home improvement. It allows users to search for properties for sale or rent, as well as estimate the value of homes through its "Zestimate" feature, which uses an algorithm to provide an estimated market value based on various factors. Zillow also offers tools for homeowners, renters, real estate agents, and property managers to list properties, advertise, and manage their real estate transactions.
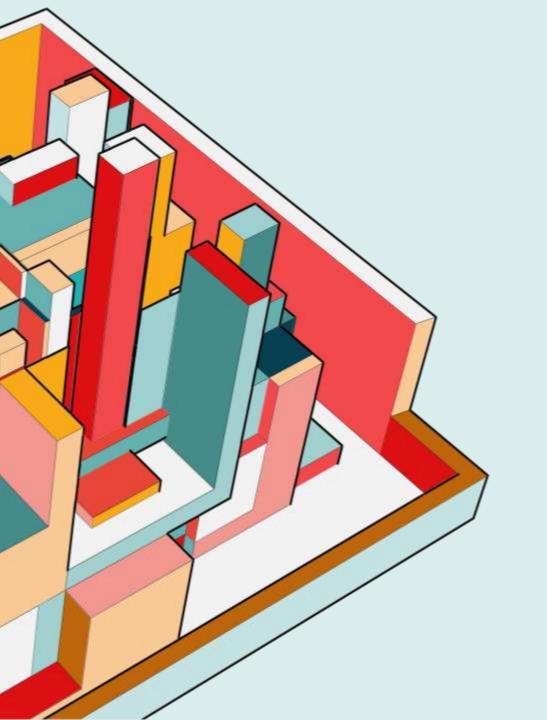
# WHY ZILLOW?

1. **Rich Dataset:** Zillow provides extensive real estate data including listings, property values, and demographic information.

2. **High-Quality Data:** Zillow's data is well-structured and reliable, simplifying analysis and modeling.

3. **Popularity:** Widely used by stakeholders, Zillow data is familiar and relevant, encouraging collaboration.

4. **Real-World Applications:** Real estate data from Zillow enables practical analysis like price prediction and trend identification.

5. **Business Insights**: Analyzing Zillow data uncovers valuable insights for investors, developers, and policymakers, creating business opportunities.

# WHY WE CHOSE LOS ANGELES?

- Los Angeles has a wide variety of property types and price ranges, from affordable homes to luxury estates. This diversity provides a rich dataset for analyzing different factors affecting housing prices and understanding various market segments.Tone inflection

- Scraping data for the entire United States can be resource-intensive in terms of computing power, time, and cost. Limiting the scope to Los Angeles may be more manageable within available resources.

- Focusing on a specific region like Los Angeles could be a pilot study to test methodologies or algorithms before scaling up to include data from other states.

# OUR MAIN GOAL

1. **Data Analysis**: Perform analysis on the collected real estate data to explore trends, patterns, and correlations between house features and prices. This phase involves examining the data to derive insights that can inform the modeling process.

2. **Model Development**: Build a predictive model that estimates house prices based on various home features. This model aims to accurately predict house prices using machine learning algorithms.

3. **Model Evaluation and Fine-Tuning**: Evaluate the performance of the predictive model using appropriate metrics and fine-tune it to improve accuracy and robustness. This iterative process involves adjusting model parameters and features based on performance results.

# STEP 1: Scraping the links of the houses

# STEP 2: Scraping the details of each house

# PREPROCESS STEPS

- Removing Duplicates

- Checking the number of missing values

- Handling Home Status

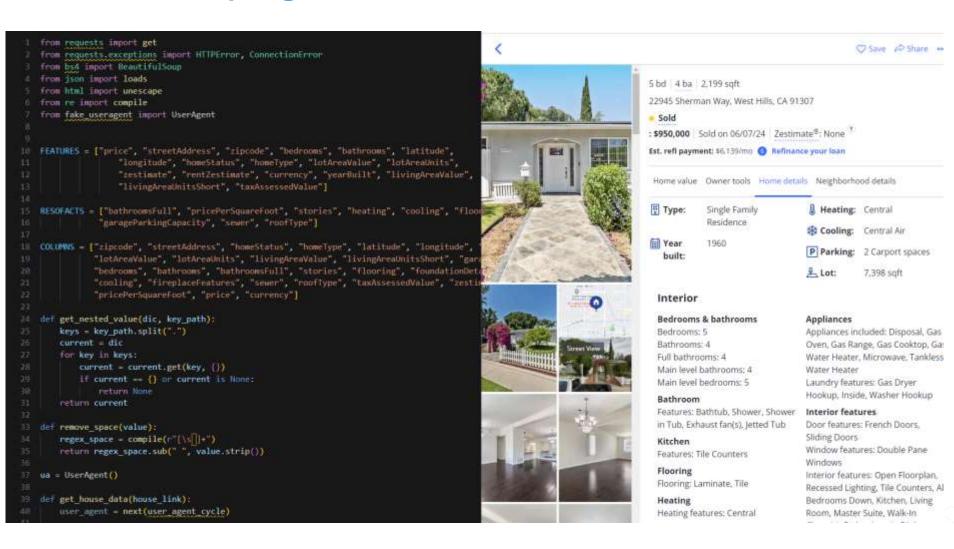- In cases where columns had a high number of null values, they were dropped from the dataset to maintain data accuracy

- We convert string columns to real arrays and merge specific values for streamlined processing.

- We merged similar types of some features like cooling together and consolidated them into broader categories

```python
1  df['cooling'] = df['cooling'].apply(lambda x: ['Central' if item in ['Central Air', 'Central Air/Evap','Central Air/Refrig','Central Forced Air'] else item for item in x])
2  df['cooling'] = df['cooling'].apply(lambda x: ['Evaporative' if item in ['Evaporative Cooling'] else item for item in x])
3  df['cooling'] = df['cooling'].apply(lambda x: ['SEER Rated' if item in ['SEER Rated 13-15','SEER Rated 16+'] else item for item in x])
4  df['cooling'] = df['cooling'].apply(lambda x: ['Wall' if item in ['Wall A/C Units','Wall Unit(s)','Wall/Window Unit(s)'] else item for item in x])
5  df['cooling'] = df['cooling'].apply(lambda x: [item for item in x if item != 'See Remarks'])
6  all_values = [item for sublist in df['cooling'] for item in sublist]
7  unique_values = set(all_values)
8  unique_values
```

- Preprocess Home Features(some examples of how we handle different features)

  - Lot area and living:
    - convert units to the square

  - Tax Assesed Value
    - Missing values in these columns were filled using a KNN imputer technique.

  - Has… Features
    - We have some columns that indicate if house has a feature or not. We will convert them to 1 or 0 and fill the rows with null values with mode and if they are too much we drop the column.

# ANALYZING DISTRIBUTIONS AND OUTLIERS

We used box plots and histograms to analyze feature distributions, identify outliers, and clean the dataset. Selecting the number of bins for histograms is crucial, and we implemented three methods:

- Sturges' Rule
- Scott's Rule
- Freedman-Diaconis Rule

# TOP 10S

```
Top 10 most expensive areas for APARTMENT based on average house prices:
+-----+---------+-----------+
|     | zipcode |   price   |
+-----+---------+-----------+
| 186 | 90291.0 | 2527557.5 |
| 127 | 90049.0 | 1350000.0 |
| 121 | 90047.0 | 850000.0  |
| 156 | 90067.0 | 773680.0  |
| 307 | 91411.0 | 699000.0  |
| 170 | 90094.0 | 685000.0  |
| 49  | 90024.0 | 659207.0  |
| 34  | 90017.0 | 208000.0  |
+-----+---------+-----------+

Top 10 most expensive areas for CONDO based on average house prices:
+-----+---------+--------------------+
|     | zipcode |       price        |
+-----+---------+--------------------+
| 157 | 90067.0 |     1807400.0      |
| 168 | 90077.0 |     1527000.0      |
| 50  | 90024.0 | 1503846.2063492064 |
| 191 | 90292.0 |     1300625.0      |
| 101 | 90041.0 |     1300000.0      |
| 128 | 90049.0 | 1255558.0232558139 |
| 124 | 90048.0 | 1239989.3043478262 |
| 249 | 91326.0 |     1212500.0      |
| 22  | 90010.0 |     1190150.0      |
| 171 | 90094.0 |     1190000.0      |
+-----+---------+--------------------+

Top 10 most expensive areas for MULTI_FAMILY based on average house prices:
+-----+---------+--------------------+
|     | zipcode |       price        |
+-----+---------+--------------------+
| 129 | 90049.0 |     4825000.0      |
| 291 | 91402.0 |     3590000.0      |
| 195 | 90293.0 |     3369000.0      |
| 54  | 90025.0 | 2948677.777777778  |
| 295 | 91403.0 |     2855000.0      |
| 320 | 91601.0 |     2562545.0      |
| 94  | 90038.0 | 2491666.6666666665 |
| 125 | 90048.0 | 2463678.5714285714 |
| 145 | 90064.0 |     2462500.0      |
| 179 | 90247.0 |     2155000.0      |
+-----+---------+--------------------+

Top 10 most expensive areas for SINGLE_FAMILY based on average house prices:
+-----+---------+--------------------+
|     | zipcode |       price        |
+-----+---------+--------------------+
| 43  | 90020.0 |     9312500.0      |
| 166 | 90069.0 | 5561361.074074074  |
| 184 | 90272.0 | 5297068.645833333  |
| 158 | 90067.0 |     4800000.0      |
| 174 | 90210.0 | 4735345.903225807  |
| 130 | 90049.0 | 4516132.431034483  |
| 51  | 90024.0 |     4514328.0      |
| 169 | 90077.0 |     3848768.0      |
| 315 | 91436.0 |     3208421.875     |
| 119 | 90046.0 | 3169094.173076923  |
+-----+---------+--------------------+

Top 10 most expensive areas for TOWNHOUSE based on average house prices:
+-----+---------+--------------------+
|     | zipcode |       price        |
+-----+---------+--------------------+
| 167 | 90069.0 |     1800000.0      |
| 190 | 90291.0 |     1767500.0      |
| 131 | 90049.0 |     1525000.0      |
| 159 | 90067.0 |     1427000.0      |
| 173 | 90094.0 |     1379500.0      |
| 185 | 90272.0 | 1324666.6666666667 |
| 23  | 90010.0 |     1310000.0      |
| 56  | 90025.0 | 1275571.4285714286 |
| 193 | 90292.0 | 1259678.5714285714 |
| 163 | 90068.0 |     1255000.0      |
+-----+---------+--------------------+
```
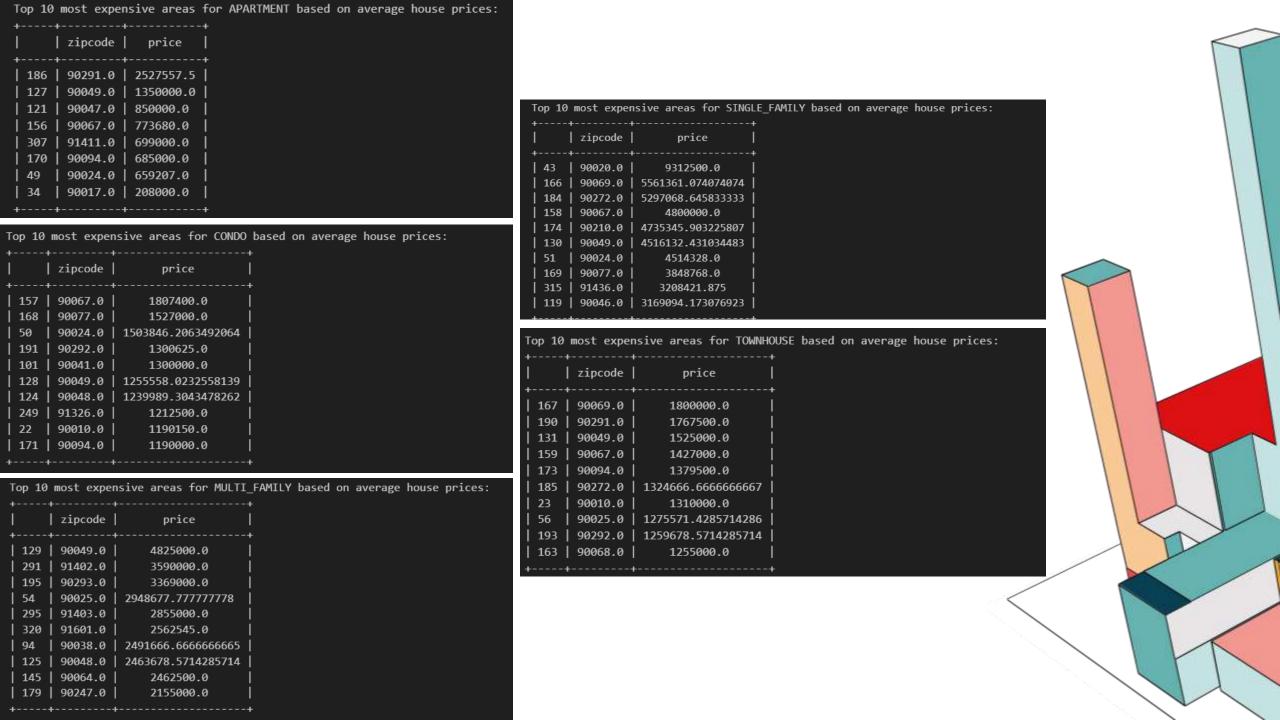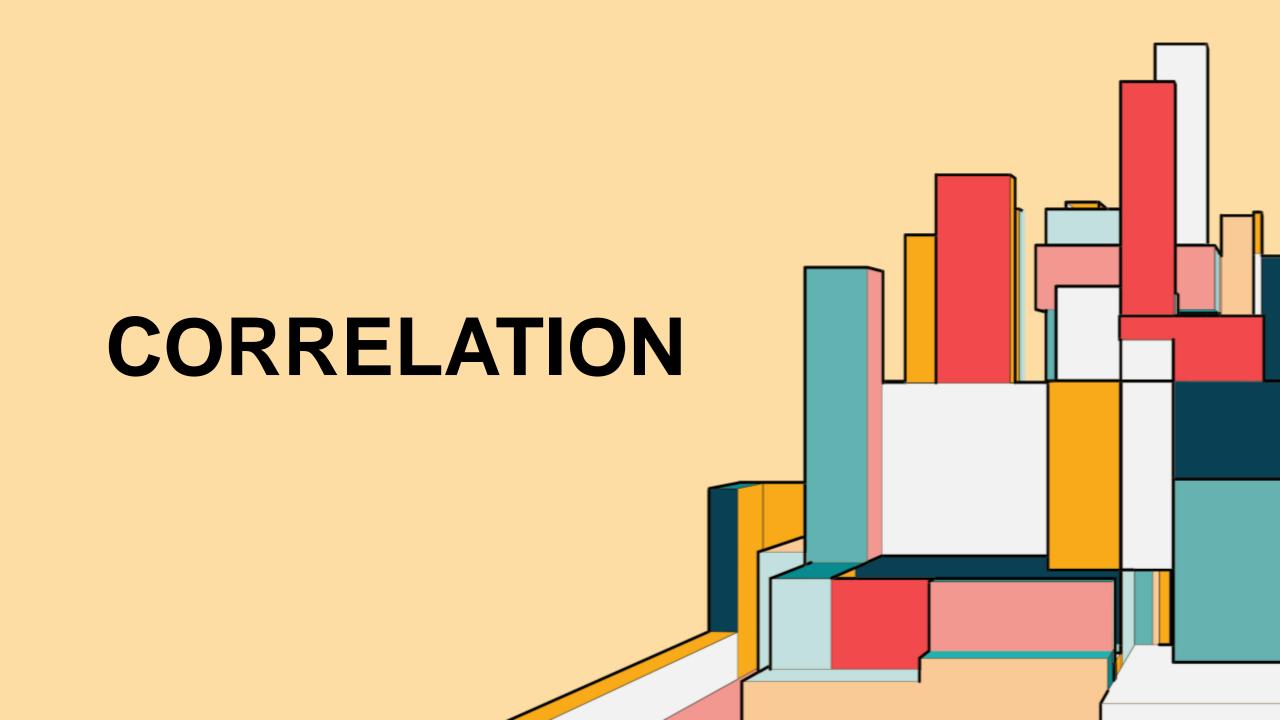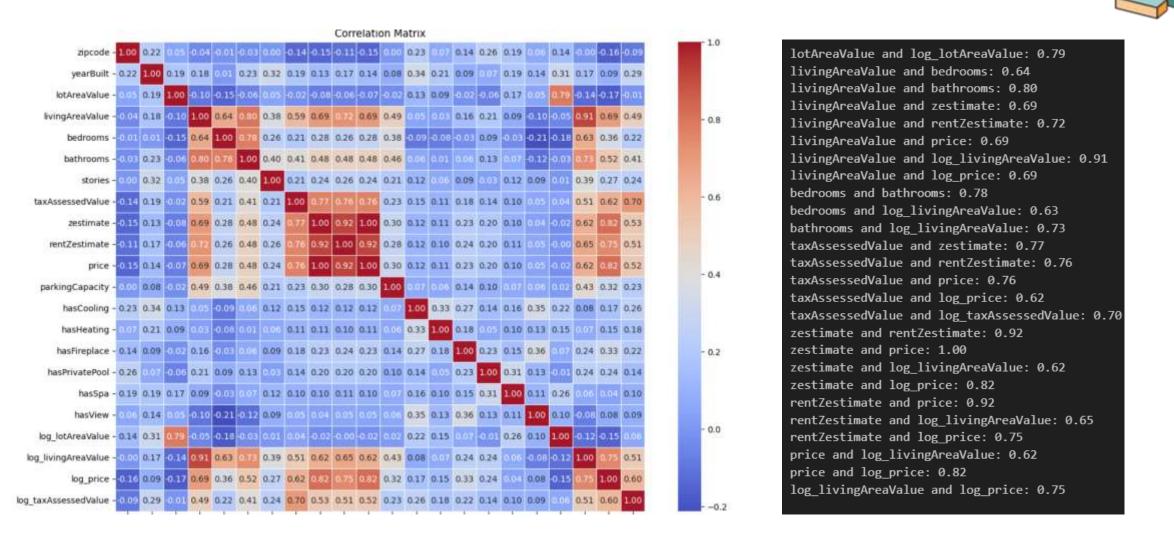
# FINDINGS FROM VISUALIZING FEATURE CHANGES BASED ON YEAR BUILT

**Here are our findings summarized:**

**- Average Bedrooms:**
   **- 1900-1940: Stabilized around 3-5 bedrooms.**
   **- 1940-1980: Increased fluctuations, peaks around 1955 and 1975.**
   **- 1980-2000: Slight downward trend, dipping close to 3 around 1995.**
   **- 2000-2020: Upward trend, reaching nearly 6 by 2020.**

**- Number of Stories:**
   **- Initially low variability and stability.**
   **- Gradual increase mid-20th century.**
   **- Significant upward trend from 1980 onwards.**

**- Cooling Systems:**
   **- Pre-1925: Almost no adoption.**
   **- 1925-1965: Significant increase, peaks around 1940 and 1955.**
   **- 1965-2000: Considerable fluctuations but overall increase.**
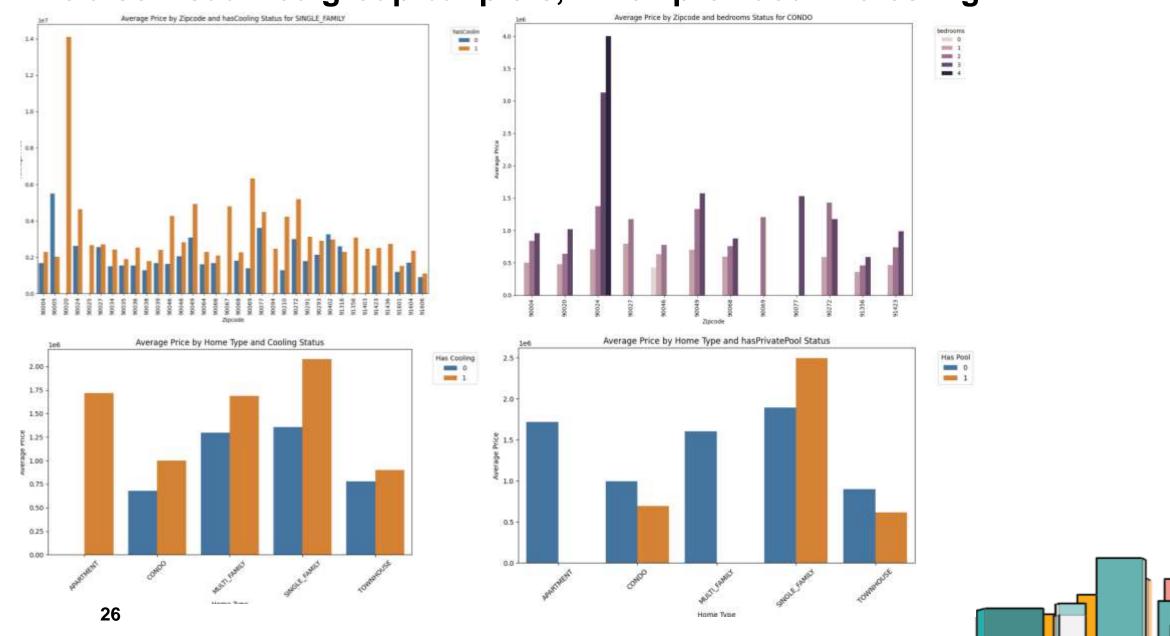   **- 2000-2020: Lower counts early on, sharp increase around 2020.**

# CORRELATION

Correlation Matrix

```
lotAreaValue and log_lotAreaValue: 0.79
livingAreaValue and bedrooms: 0.64
livingAreaValue and bathrooms: 0.80
livingAreaValue and zestimate: 0.69
livingAreaValue and rentZestimate: 0.72
livingAreaValue and price: 0.69
livingAreaValue and log_livingAreaValue: 0.91
livingAreaValue and log_price: 0.69
bedrooms and bathrooms: 0.78
bedrooms and log_livingAreaValue: 0.63
bathrooms and log_livingAreaValue: 0.73
taxAssessedValue and zestimate: 0.77
taxAssessedValue and rentZestimate: 0.76
taxAssessedValue and price: 0.76
taxAssessedValue and log_price: 0.62
taxAssessedValue and log_taxAssessedValue: 0.70
zestimate and rentZestimate: 0.92
zestimate and price: 1.00
zestimate and log_livingAreaValue: 0.62
zestimate and log_price: 0.82
rentZestimate and price: 0.92
rentZestimate and log_livingAreaValue: 0.65
rentZestimate and log_price: 0.75
price and log_livingAreaValue: 0.62
price and log_price: 0.82
log_livingAreaValue and log_price: 0.75
```

# GROUP BAR PLOTS

# We also visualized group bar plots, which provided interesting



26

# Examples of Insights from Group Bar Plot Analysis

**Single Family Homes:**

1.**Desirable Views**: Highly desirable, premium pricing.

2.**Large Properties**: Better utilization of views.

3.**Luxury Status**: Associated with prestige and luxury.

**Condos:**

1.**Minimal View Impact**: Shared amenities, urban settings.

2.**Urban Proximity**: Convenience trumps views.

3.**Cost Focus**: Buyers prioritize affordability and location.

**Townhouses/Apartments:**

1.**Moderate View Value**: Higher floors, less exclusivity.

2.**Urban/Suburban**: Less impact in urban, some value in suburban.

3.**Convenience First**: Buyers prefer amenities and maintenance over views.

**Insights from High-Priced Zip Codes:**

•**Single Family Homes**: Views significantly increase prices due to desirability and exclusivity.

•**Townhouses**: Higher prices linked to scenic locations.

•**Condos and Multifamily Homes**: Views have minimal impact; shared amenities and property features are more important.

This highlights the need to consider property type and key value factors when assessing the impact of views on home prices in high-demand areas.

**Parking Impact Insights:**

**Single Family Homes:**

1.**Convenience**: Ample parking boosts demand and prices.

2.**Value**: Adds significant property value.

**Condos and Multifamily Homes:**

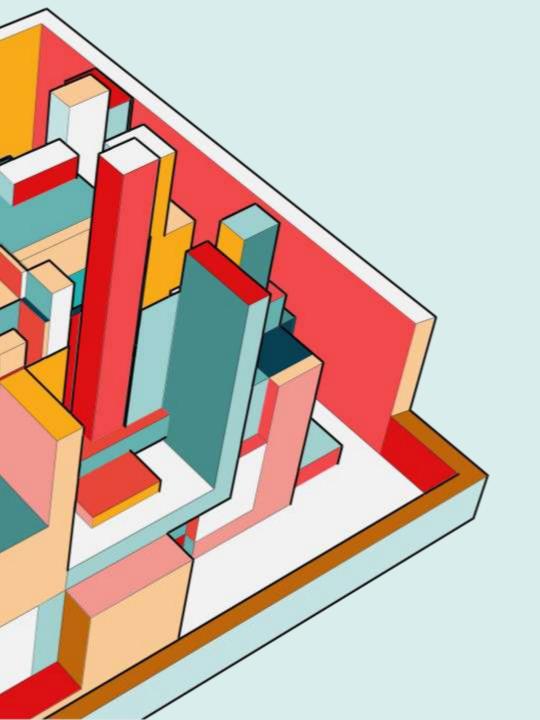1.**Urban Necessity**: Parking crucial in urban areas, affecting prices.

2.**Demand**: High parking demand impacts design and pricing.

**Townhouses and Apartments:**

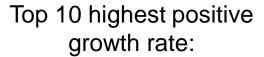1.**Transport Access**: Less reliant on parking due to public transit.

2.**Shared Facilities**: Shared parking reduces focus on individual parking spaces.

GROWTH RATE BASED ON SOLD-HISTORY

We parse the sold history of homes, calculate annual growth rates for each property's price, and then compute the average growth rate by zipcode.

Top 10 highest positive growth rate:

```
90047 0.172034
90065 0.158968
90062 0.126168
90037 0.096886
90064 0.093413
90230 0.093221
91605 0.089674
90008 0.085783
90061 0.080442
90011 0.076536
```

Top 10 highest negative growth rate:

```
90019  -1.020973
90501  -1.255977
90067  -1.265953
91604  -1.282510
90272  -1.730388
90017  -3.617703
91601  -4.139532
90028  -5.974681
90013  -7.372272
90015 -13.803312
```

# A brief exploration into these neighborhoods🚀:

**POSOTIVE GROWTH RATE:**

- **Urban Revitalization**:
  - **90047, 90062**: South LA sees revitalization efforts and new housing developments.
  - **90037**: Near USC, demand from students and faculty boosts property values.
- **Improved Accessibility**:
  - **90064**: Metro Expo Line expansion enhances accessibility in Rancho Park.
  - **90230**: Culver City's redevelopment attracts residential and commercial interest.
- **Employment and Amenities**:
  - **91605**: NoHo Arts District's cultural attractions appeal to young professionals.
  - **90011**: Improving amenities draw residents to South LA near Central-Alameda.
- **Community and Culture**:
  - **90008**: Baldwin Hills' cultural significance attracts residents.
  - **90065**: Scenic views and community vibe in Glassell Park appeal to families.
- **Affordability**:
  - **90061**: Watts area offers relatively affordable housing options.
  - **90047**: South LA provides affordable living compared to other parts of LA.

**NEGATIVE GROWTH RATE:**

1. **90019 (Mid-City)**: Limited development space and higher crime rates.
2. **90501 (Torrance)**: Stable area with limited inventory and slower growth.
3. **90067 (Century City)**: High property prices and limited residential options.
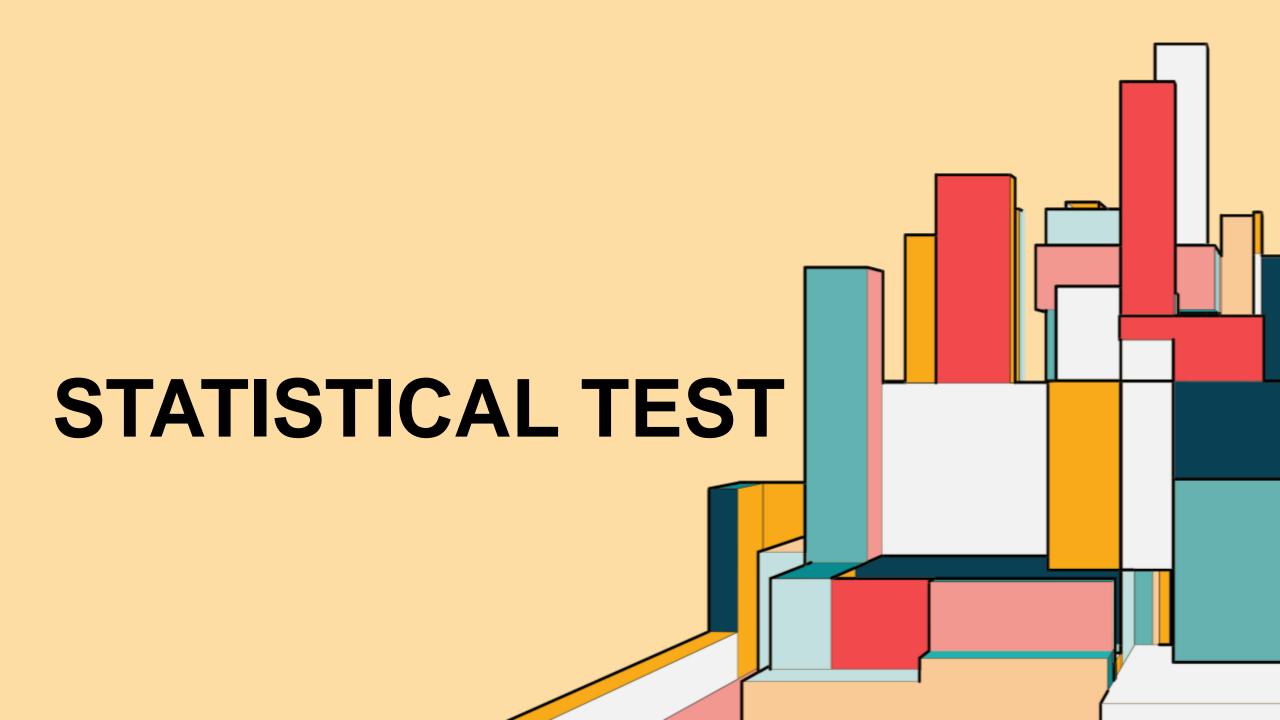4. **91604 (Studio City)**: High prices and limited inventory, leading to competition.
5. **90272 (Pacific Palisades)**: Strict zoning and limited land availability.
6. **90017 (Downtown LA - Historic Core)**: Homelessness and urban blight issues.
7. **91601 (North Hollywood - Arts District)**: Industrial past and crime concerns.
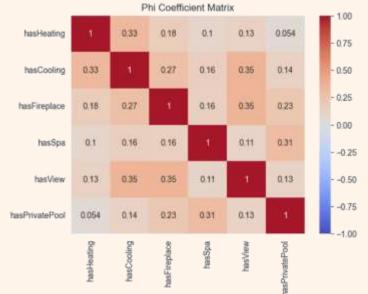8. **90028 (Hollywood)**: High crime rates, traffic congestion, and homelessness.
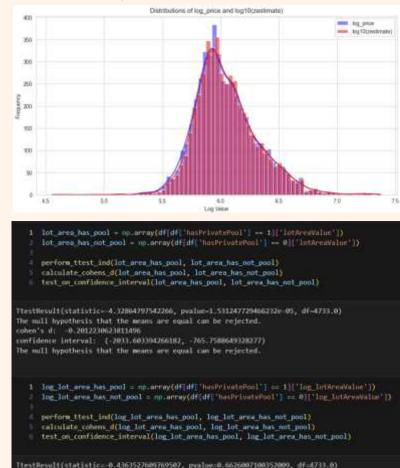9. **90013 (Downtown LA - Fashion District)**: Industrial area with limited amenities.
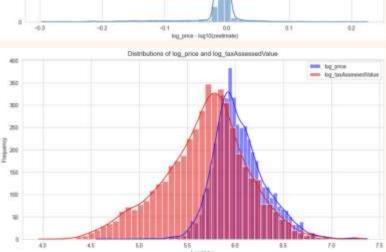10. **90015 (Downtown LA - South Park)**: Dense urban area with few residential amenities.

# STATISTICAL TEST

We will use the following statistical tests to analyze the features and their relationships:

- **Paired T-Test**

- **Independent T-Test**

- **Confidence Interval**

- **Chi-Square Test**

- **Cramér's V and Phi Coefficient**



Distributions of log_price and log10(zestimate)



Difference between log_price and log10(zestimate)



Phi Coefficient Matrix

```
1  lot_area_has_pool = np.array(df[df['hasPrivatePool'] == 1]['lotAreaValue'])
2  lot_area_has_not_pool = np.array(df[df['hasPrivatePool'] == 0]['lotAreaValue'])
3
4  perform_ttest_ind(lot_area_has_pool, lot_area_has_not_pool)
5  calculate_cohens_d(lot_area_has_pool, lot_area_has_not_pool)
6  test_on_confidence_interval(lot_area_has_pool, lot_area_has_not_pool)
```

```
TtestResult(statistic=-4.32864797542266, pvalue=1.531247729466232e-05, df=4733.0)
The null hypothesis that the means are equal can be rejected.
cohen's d:  -0.2012230623011496
confidence interval:  (-2033.603394266182, -765.7588649328277)
The null hypothesis that the means are equal can be rejected.
```

```
1  log_lot_area_has_pool = np.array(df[df['hasPrivatePool'] == 1]['log_lotAreaValue'])
2  log_lot_area_has_not_pool = np.array(df[df['hasPrivatePool'] == 0]['log_lotAreaValue'])
3
4  perform_ttest_ind(log_lot_area_has_pool, log_lot_area_has_not_pool)
5  calculate_cohens_d(log_lot_area_has_pool, log_lot_area_has_not_pool)
6  test_on_confidence_interval(log_lot_area_has_pool, log_lot_area_has_not_pool)
```

```
TtestResult(statistic=-0.4363527609769507, pvalue=0.6626007100352009, df=4733.0)
The null hypothesis that the means are equal cannot be rejected.
cohen's d:  -0.02028440948501723
confidence interval:  (-0.05464008881540952, 0.03476029458423553)
The null hypothesis that the means are equal cannot be rejected.
```

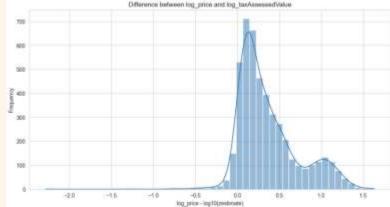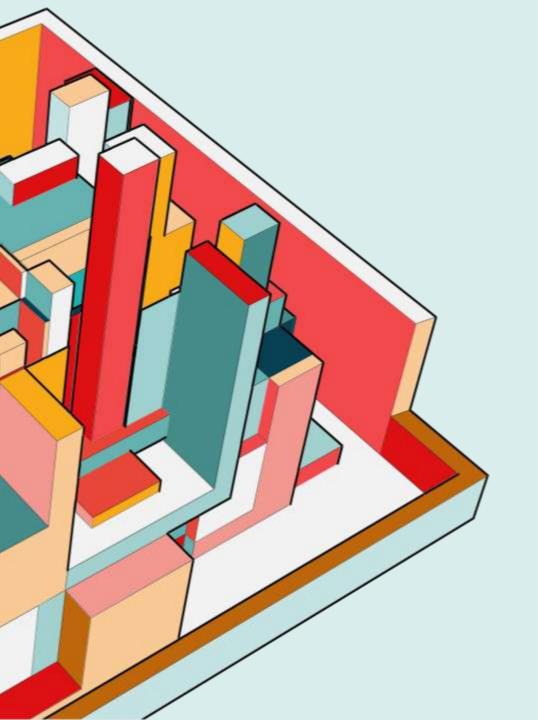Distributions of log_price and log_taxAssessedValue



Difference between log_price and log_taxAssessedValue



```
1  perfrom_chi2_test(df_filtered_hometypes["homeType"], df_filtered_hometypes["hasPrivatePool"])
```

```
chi2:  249.5129117256268
p:  8.340199418475354e-54
The null hypothesis that the variables are independent can be rejected.
crammers_v:  0.22875085839539938
```

NEXT STEP....

Our future work involves developing a predictive model to estimate house prices using various features, such as location, size, number of bedrooms, and other relevant attributes. This model aims to improve accuracy in price predictions by leveraging the detailed data we've gathered and analyzed.

THANK YOU !