



# Exploratory Data Analysis (EDA) & Data Cleaning for House Pricing Dataset

PREPARED BY VAISHALI PATELIA

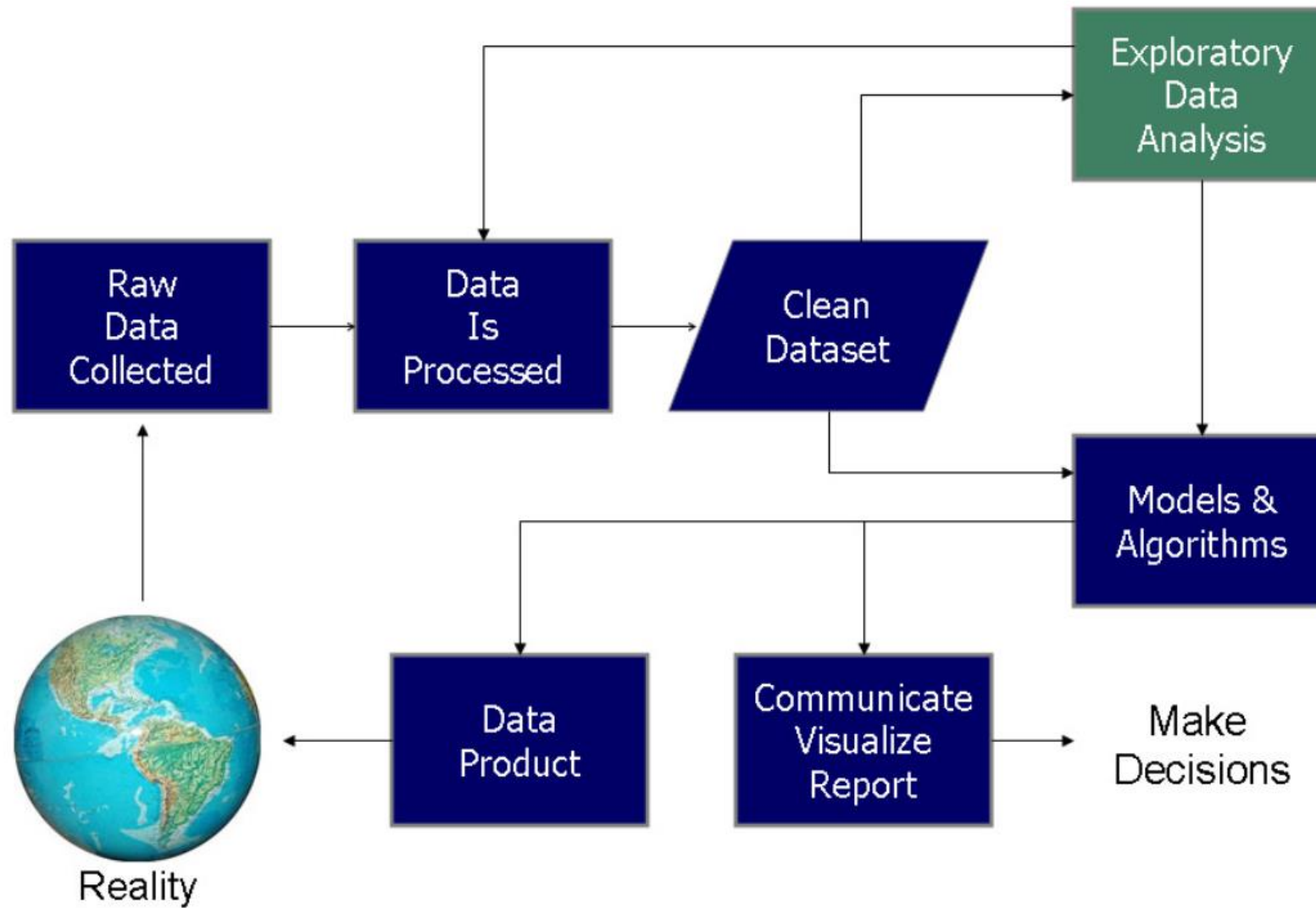
# Outline

- ▶ Introduction
- ▶ Purpose of the Project
- ▶ Dataset Overview
- ▶ Python packages for EDA
- ▶ Data Cleaning
- ▶ Data Visualization
- ▶ Conclusion
- ▶ Recommendation
- ▶ References
- ▶ Appendices



# Data Science Process

3



Source: <https://datasilk.com/data-analysis/>

# Purpose of the Project

**Goal:** Clean and prepare a housing dataset for the modeling team.

**Tasks Performed:**

1. Data Loading & Initial Analysis
2. Handling Missing Values
3. Data Type Conversion
4. Outlier Detection & Removal
5. Data Visualization & Insights

# Dataset Overview

- ▶ Total Rows & Columns: Displayed from `df.shape` (5000, 16)
- ▶ Feature Summary: Key columns (e.g., `sold_price`, `sqrt_ft`, `bedrooms`, `lot_acres`)
- ▶ Missing Values: Shown using `df.info()` (`lot_acres`, `bathrooms`, `sqrt_ft`, `garage`, `HOA`)
- ▶ Data Types: Numeric & Categorical columns (`bathrooms`, `sqrt_ft`, `garage`, `fireplaces`, `HOA`)

`df.head()`

	MLS	sold_price	zipcode	longitude	latitude	lot_acres	taxes	year_built	bedrooms	bathrooms	sqrt_ft	garage	kitchen_features	fireplaces	floor_cov
0	21530491	5300000.0	85637	-110.378200	31.356362	2154.00	5272.00	1941	13	10.0	10500.0	0.0	Dishwasher, Freezer, Refrigerator, Oven	6.0	Mexican V
1	21529082	4200000.0	85646	-111.045371	31.594213	1707.00	10422.36	1997	2	2.0	7300.0	0.0	Dishwasher, Garbage Disposal	5.0	Natural S C
2	3054672	4200000.0	85646	-111.040707	31.594844	1707.00	10482.00	1997	2	3.0	NaN	NaN	Dishwasher, Garbage Disposal, Refrigerator	5.0	Natural S Other:
3	21919321	4500000.0	85646	-111.035925	31.645878	636.67	8418.58	1930	7	5.0	9019.0	4.0	Dishwasher, Double Sink, Pantry: Butler, Refri...	4.0	Ceramic Lami V
4	21306357	3411450.0	85750	-110.813768	32.285162	3.21	15393.00	1995	4	6.0	6396.0	3.0	Dishwasher, Garbage Disposal, Refrigerator, Mi...	5.0	Ce Con

# Python packages for EDA



# Data Cleaning

- ▶ **Handled Missing Values** : Used median imputation for numerical data & 0 for categorical values
- ▶ **Outlier Detection & Removal** : Applied Interquartile Range (IQR) method
- ▶ **Fixed Data Types** : Converted numerical values stored as strings
- ▶ **Removed Duplicates** : Ensured data consistency
- ▶ **Standardized Column Names** : For easy reference in modeling

## Summary Statistics After Cleaning:

	MLS	sold_price	zipcode	longitude	latitude \
count	3.733000e+03	3.733000e+03	3733.000000	3733.000000	3733.000000
mean	2.135377e+07	6.837639e+05	85726.135548	-110.916177	32.324772
std	2.101180e+06	1.362301e+05	32.688514	0.092688	0.134497
min	3.042851e+06	3.750000e+05	85118.000000	-111.430863	31.458609
25%	2.140831e+07	5.750000e+05	85718.000000	-110.975535	32.285978
50%	2.161784e+07	6.500000e+05	85737.000000	-110.922752	32.319066
75%	2.180678e+07	7.500000e+05	85750.000000	-110.861144	32.396889
max	2.192856e+07	1.185000e+06	85935.000000	-109.861617	34.314889

	lot_acres	taxes	year_built	bedrooms	bathrooms \
count	3733.000000	3733.000000	3733.000000	3733.000000	3729.000000
mean	1.06910	6006.996552	1992.924993	3.850522	3.600697
std	0.84015	2092.235652	49.317624	0.826379	0.968307
min	0.00000	459.530000	0.000000	2.000000	2.000000
25%	0.51000	4708.000000	1987.000000	3.000000	3.000000
50%	0.87000	5945.000000	1999.000000	4.000000	4.000000
75%	1.20000	7272.620000	2005.000000	4.000000	4.000000
max	3.50000	11809.000000	2019.000000	18.000000	36.000000

	sqrt_ft	garage	fireplaces	HOA
count	3733.000000	3733.000000	3714.000000	3733.000000
mean	3423.058666	2.737209	1.721055	63.131372
std	617.383578	0.918898	1.000271	66.155298
min	1780.000000	0.000000	0.000000	0.000000
25%	2998.000000	2.000000	1.000000	3.000000
50%	3401.000000	3.000000	2.000000	44.000000
75%	3811.000000	3.000000	2.000000	100.000000
max	5125.000000	12.000000	8.000000	263.000000

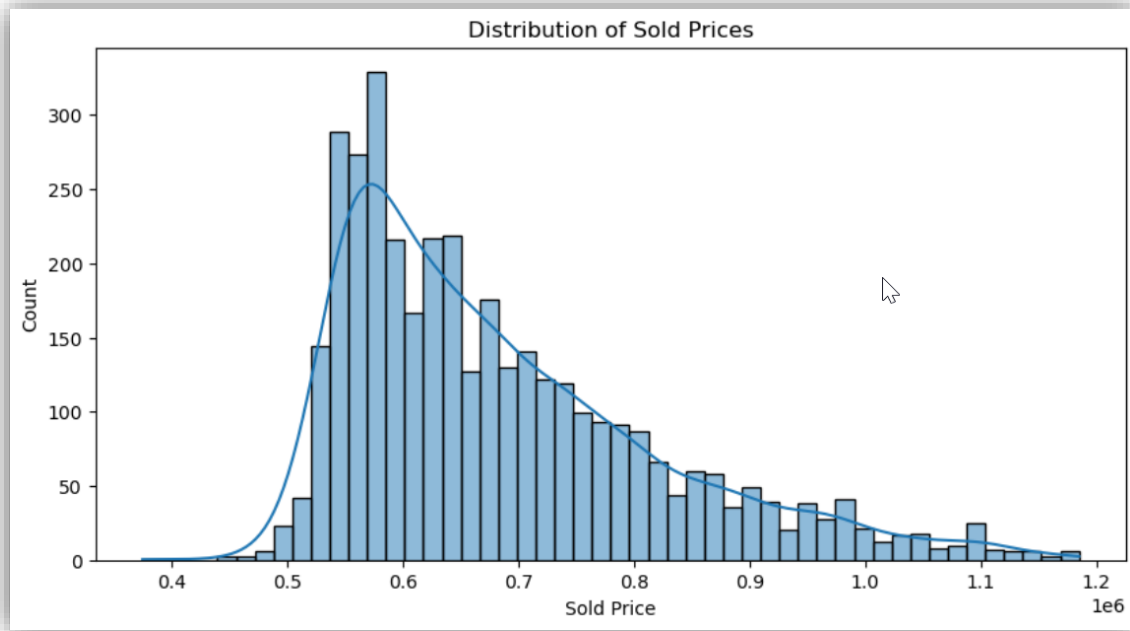
# Data Visualization

- ▶ Helps detect patterns, trends, and anomalies
- ▶ Makes data easier to interpret and communicate

## **Visuals Used in this Project:**

- ▶ Histograms & Boxplots : Understanding distributions & outliers
- ▶ Correlation Heatmap : Finding relationships between variables
- ▶ Scatterplots : Checking trends between house size & price
- ▶ Bar Chart : Compare categorical data distributions (e.g., house type, location)
- ▶ Pair Plot : Shows relationships between multiple numerical variables



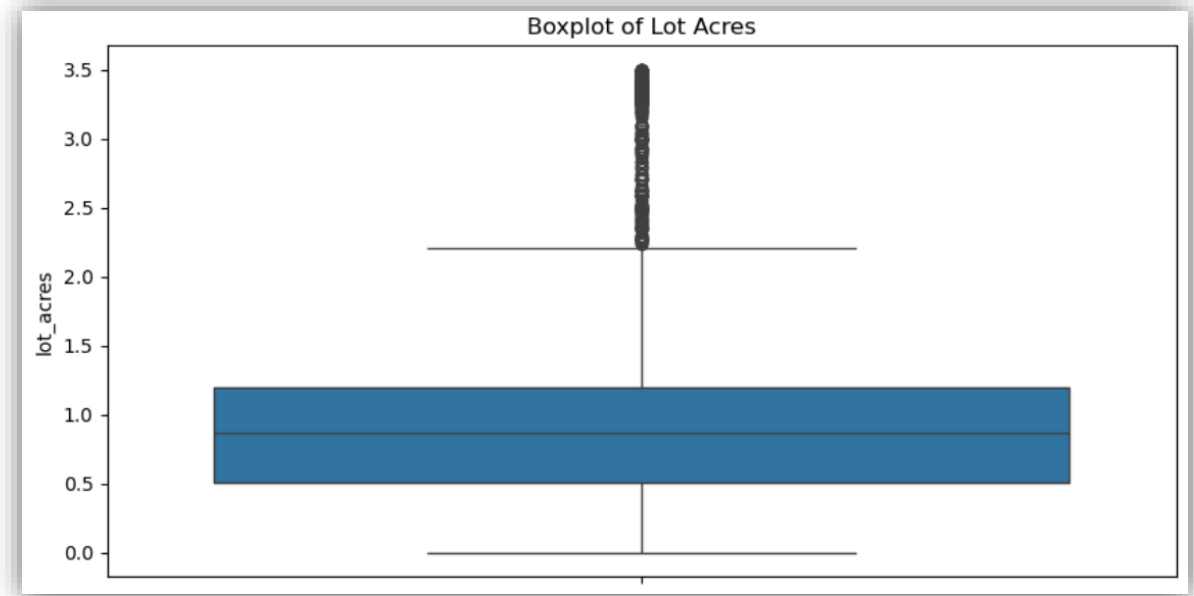


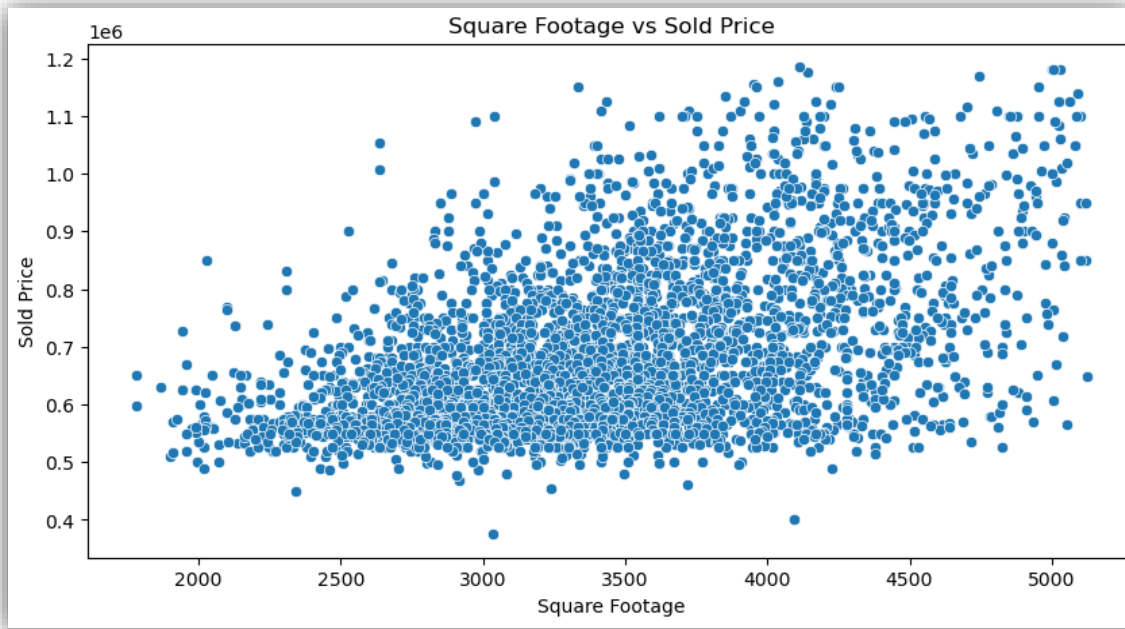
## Histogram of Sold Prices

House prices were right-skewed, meaning some high-value properties affected the average.

## Boxplot of Lot Acres

Shows data distribution & outliers using quartiles.



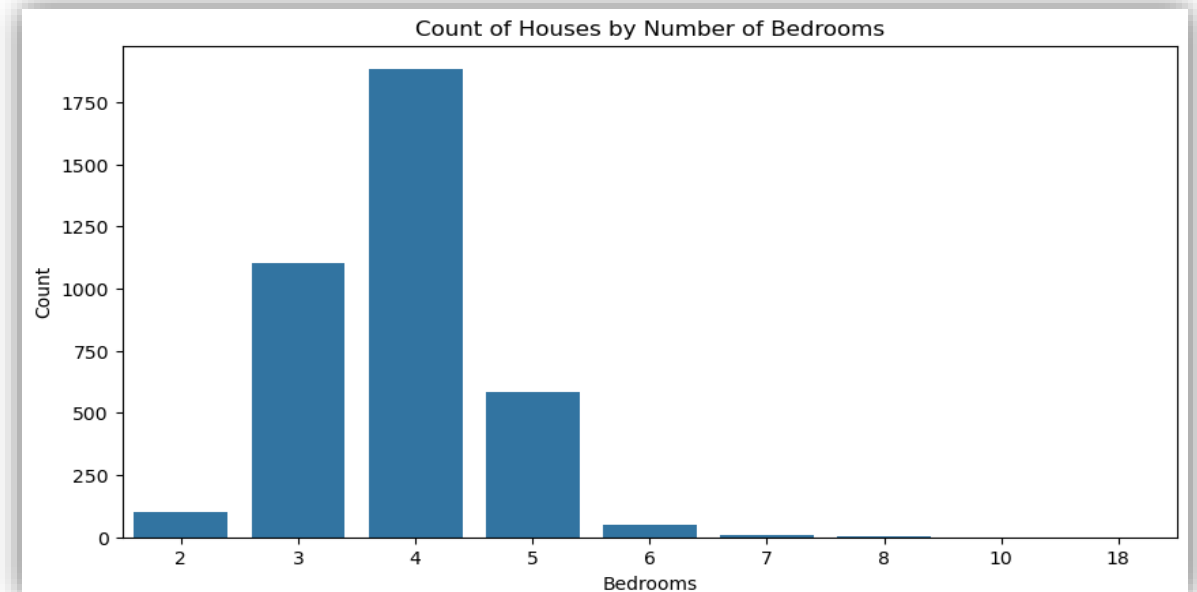


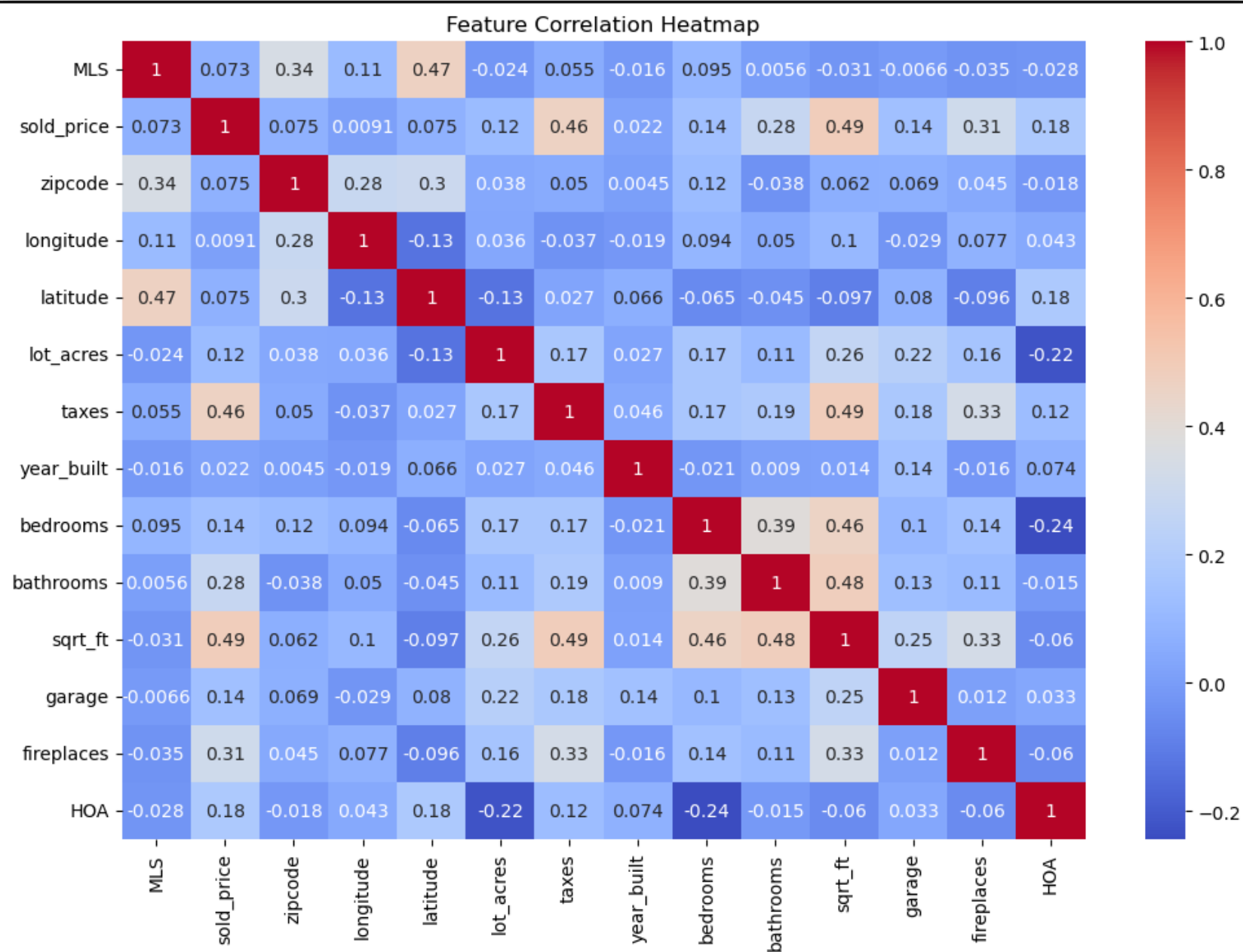
## Scatter plot of Square Footage vs Price

Used to visualize the relationship between house size (square footage) and sold price. This helps identify whether larger homes tend to have higher prices.

## Count Plot of Number of Bedrooms

Used to visualize the distribution of houses based on the number of bedrooms. This helps identify the most common house types in the dataset.





## Correlation Heatmap

Helps identify feature relationships to select relevant predictors for modeling.

Square footage, number of rooms, and overall quality had the highest correlation with house prices.

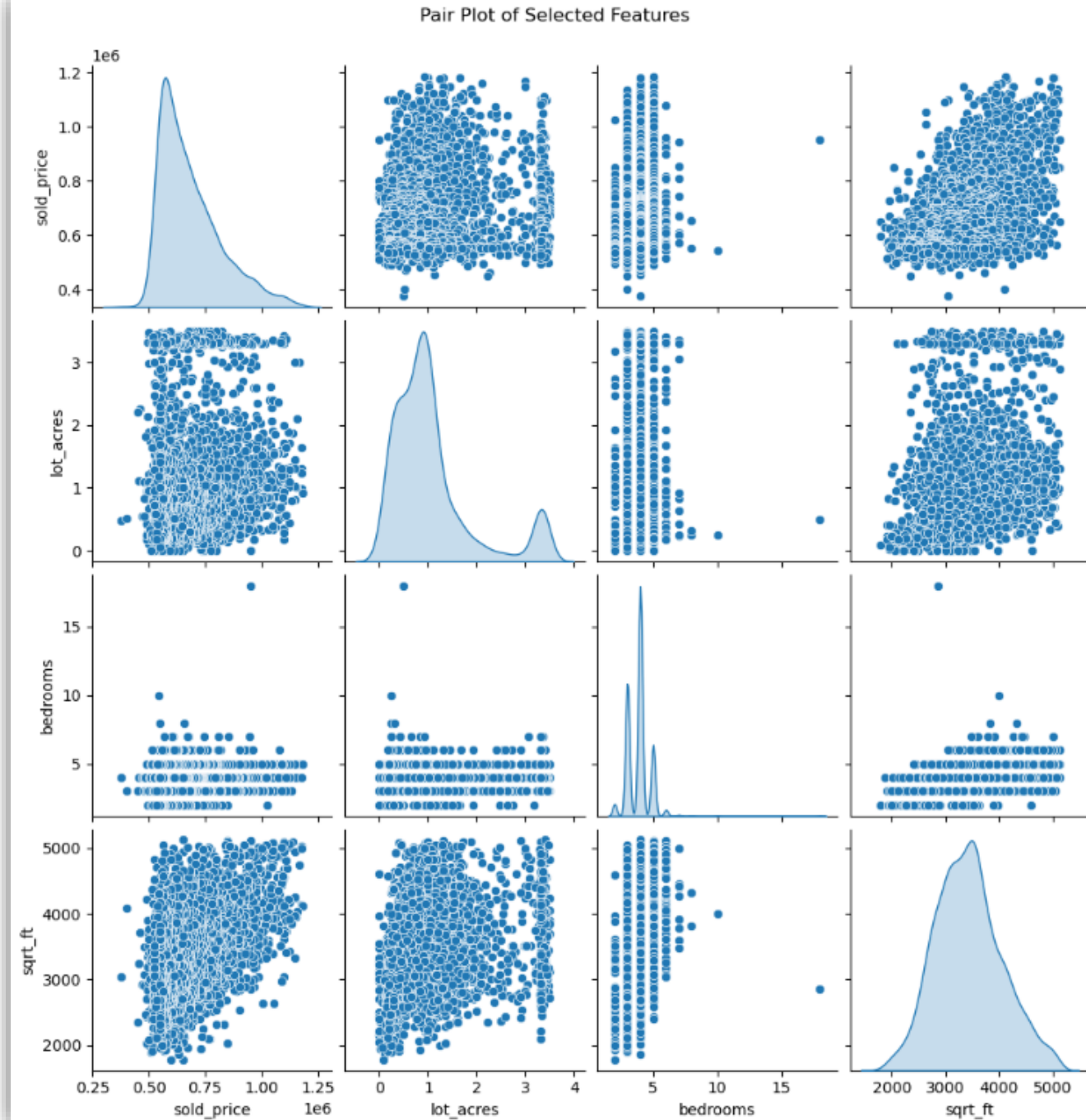
Lot size had weak correlation, indicating it might not be a strong predictor.

## PairPlot of Selected Features

House price vs. square footage showed a strong positive correlation (larger homes tend to be more expensive).

Year built vs. price suggested that newer homes generally have higher prices.

Outliers were visible in some features like lot size.



# Conclusion

- ▶ Missing values handled using median imputation
- ▶ Outliers removed using IQR method to improve model accuracy
- ▶ Right-skewed price distribution, suggesting need for log transformation
- ▶ Strong correlation between house size, number of rooms, and price
- ▶ Categorical features (e.g., location, house type) significantly impact price

# Recommendation

- ▶ Feature Engineering for model improvement
- ▶ Train machine learning models on cleaned data
- ▶ Fine-tune models for better predictions

# References

- ▶ Han, J., Kamber, M., & Pei, J. (2011). *Data Mining: Concepts and Techniques*. Elsevier.
- ▶ Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- ▶ Hawkins, D. M. (1980). *Identification of Outliers*. Springer.
- ▶ García, S., Luengo, J., & Herrera, F. (2015). *Data Preprocessing in Data Mining*. Springer.
- ▶ Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow*. O'Reilly Media.
- ▶ Pedregosa, F., et al. (2011). *Scikit-learn: Machine Learning in Python*. Journal of Machine Learning Research.
- ▶ Seabold, S., & Perktold, J. (2010). *Statsmodels: Econometric and Statistical Modeling with Python*.

# Appendices

## What is EDA?

EDA is the **process of analyzing and summarizing datasets** to uncover patterns, detect anomalies, and gain insights before building models.

### Key Objectives:

- ▶ Understand data structure and distributions
- ▶ Identify missing values and inconsistencies
- ▶ Detect and handle outliers
- ▶ Explore correlations between variables
- ▶ Generate visualizations to interpret trends

### Techniques Used in EDA:

- ▶ **Descriptive Statistics:** Mean, median, standard deviation
- ▶ **Data Visualization:** Histograms, Boxplots, Scatterplots, Heatmaps
- ▶ **Outlier Detection:** Interquartile Range (IQR), Z-score



# Handled Missing Values

## Mean

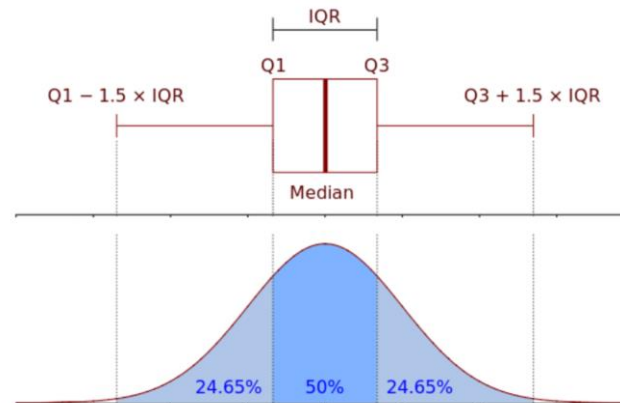
The average of all values in a dataset

Best for normally distributed data

## Median

The middle value when data is ordered

Better for skewed data or when outliers are present



## Standard Deviation (SD)

Measure of spread or dispersion of data from the mean

Lower SD = data clustered closer to the mean

Higher SD = data more spread out

Quantifies variability and consistency in data

# Outlier Detection

## Interquartile Range (IQR)

- ▶ IQR is the range between the first quartile (Q1) and third quartile (Q3)
- ▶  $IQR = Q3 - Q1$
- ▶ Outlier identification:
  - ▶ Lower bound:  $Q1 - 1.5 * IQR$
  - ▶ Upper bound:  $Q3 + 1.5 * IQR$
  - ▶ Any data point outside these bounds is considered an outlier

## Z-score

- ▶ Measures how many standard deviations a data point is from the mean
- ▶  $Z = (X - \mu) / \sigma$   
Where X is the data point,  $\mu$  is the mean, and  $\sigma$  is the standard deviation
- ▶ Outlier identification:
  - ▶ Typically,  $|Z| > 3$  is considered an outlier
  - ▶ This threshold can be adjusted based on the specific needs of the analysis



Thank You