# PREDICTIVE ANALYSIS IN PROPERTY VALUATION

A King County Case Study

# TABLE OF CONTENTS
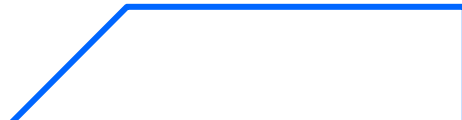
# INTRODUCTION

- Real estate agencies guide homeowners in crucial decisions like buying, selling, or renting properties.

- The project empowers agencies with a regression-based model predicting property value increases based on features like bedrooms, year built, floors, living space, condition, and location.

- The model assists in pricing strategies, market analysis, and property inspections, maximizing return on investment for clients.
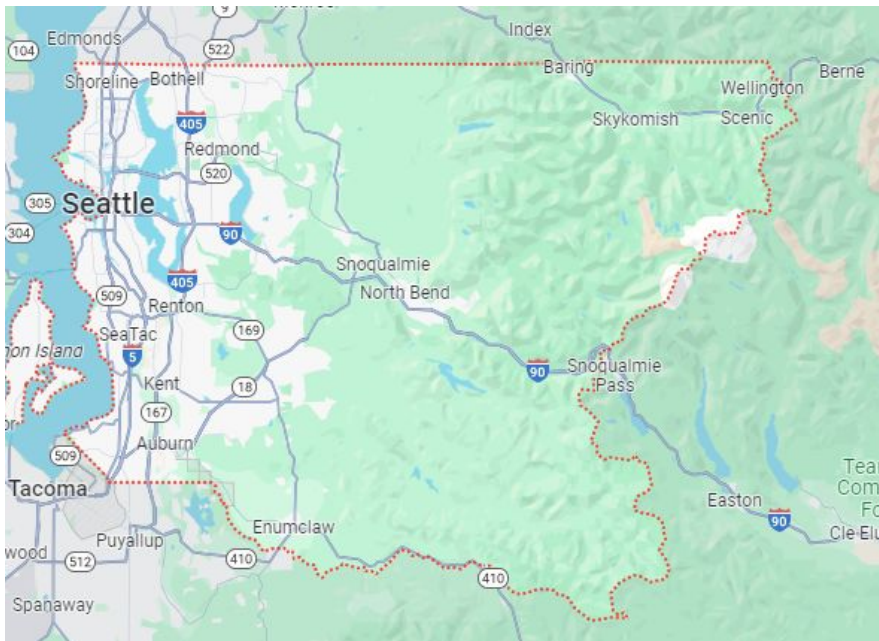
# BACKGROUND OF THE DATA

- The housing market in King County, Washington, has a rich history, marked by significant growth and changes.

- Situated in the northwestern U.S., King County includes the major economic and cultural hub of Seattle.

- King County's real estate market is recognized for its competitiveness, with diverse neighborhoods providing a range of housing options from historic homes to modern developments.

# DATA SOURCE

- King County
- Located in the US state of Washington
- 21,597 listings from May 2014 to May 2015

# PROBLEM STATEMENT

Real estate agencies struggle to guide clients on pricing, market trends, and property inspections, influenced by complex factors. The problem centers on providing agencies a precise solution for informed client recommendations.

The project aims to develop a regression model, leveraging the King County House Sales dataset, to empower agencies with evidence-based insights and optimize property value.

# TARGET AUDIENCE

- Real estate professionals seeking insights on pricing, market trends, and property inspections.
- Data scientists interested in developing regression models with the King County House Sales dataset.
- Real estate industry professionals aiming to enhance advisory capabilities.
- Stakeholders in the real estate landscape looking to optimize property value.

# OBJECTIVES

To develop a predictive regression model that assists real estate agencies in advising clients on house prices

To identify key factors influencing house prices in King County, California, to provide valuable insights for precise pricing strategies.

To analyze model performance using metrics such as mean squared error, R-squared values, and residual analysis to gauge the model's effectiveness.
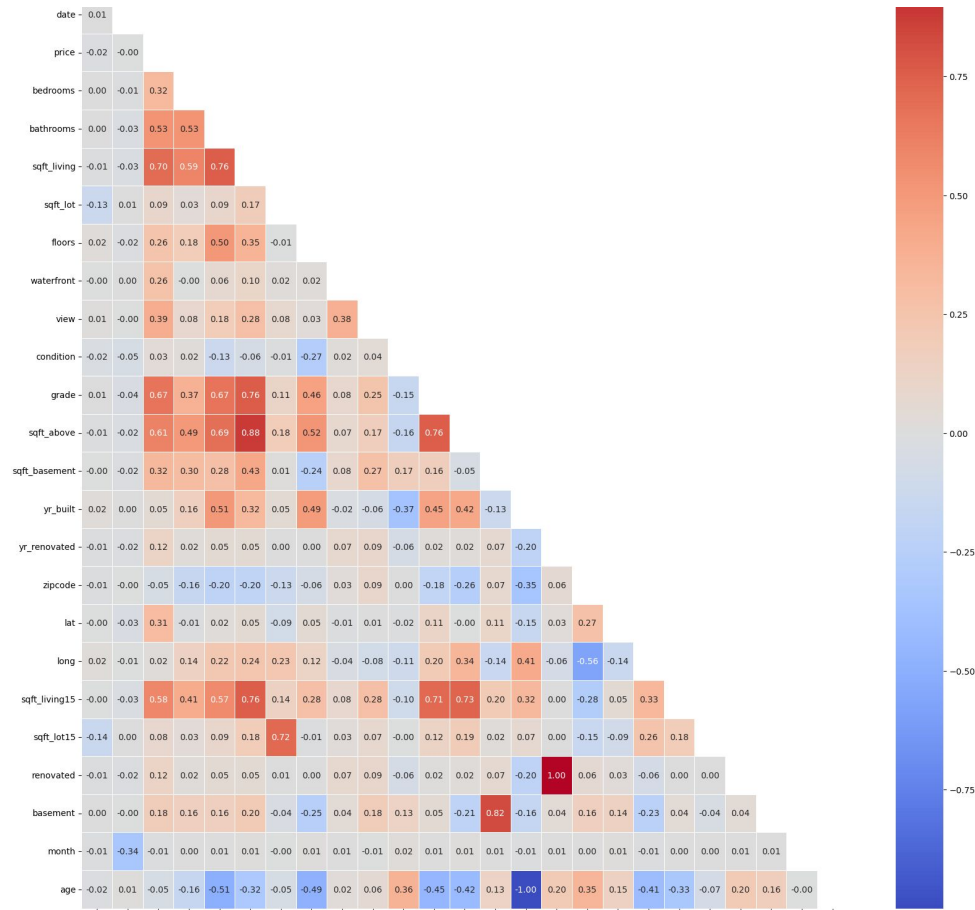
To provide actionable recommendations to the Real Estate Agency for improving profitability and market presence, leveraging insights from the model.
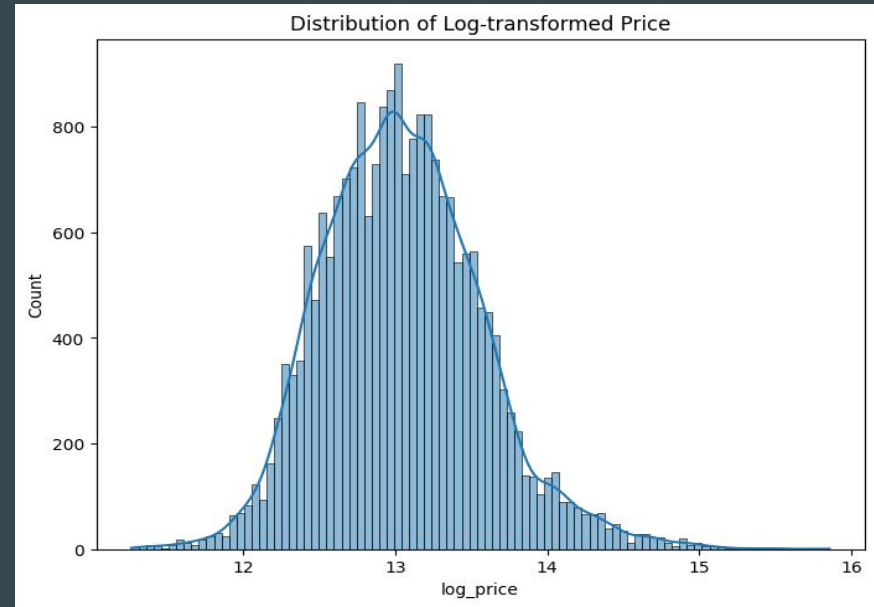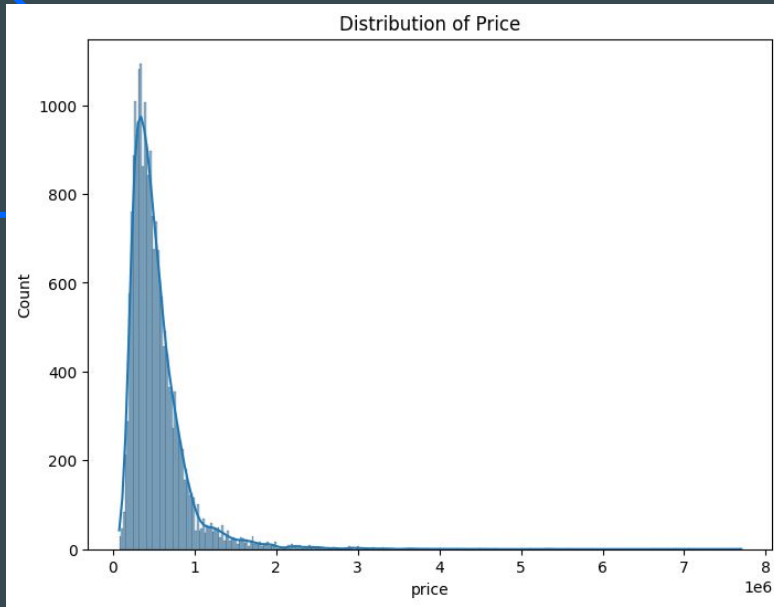
# Statistical Analysis

# Triangular Correlation Matrix

Analysis of correlation between the variables
For instance,
Year Built' and 'Age'
show a perfect negative correlation (-1.00), indicating as the property's construction year increases, its age decreases linearly.

Distribution of Price / Distribution of Log-transformed Price

The log transformation successfully normalized the positively skewed price distribution, addressing extreme values and enhancing the data's suitability for more reliable statistical analysis.

# Hypothesis Testing

## Null

There is no statistically significant relationship between the selected features and housing prices.

## Alternative

The selected features have a statistically significant relationship with housing prices.

# FINDINGS

## From hypothesis testing

## ANOVA TEST

-The null hypothesis that there is no statistically significant relationship between the selected features and house price was rejected.

-Most of the features were found to be statistically significant apart from a few which include longitude,square footage of the lot a nd square footage of the land lots of the nearest 15 neighbors.
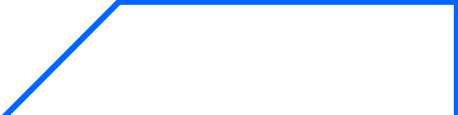
DATA MODELLING

# Models

## Baseline Model

A simple linear regression that considers a single predictor variable to estimate house price.Moderate level of prediction accuracy
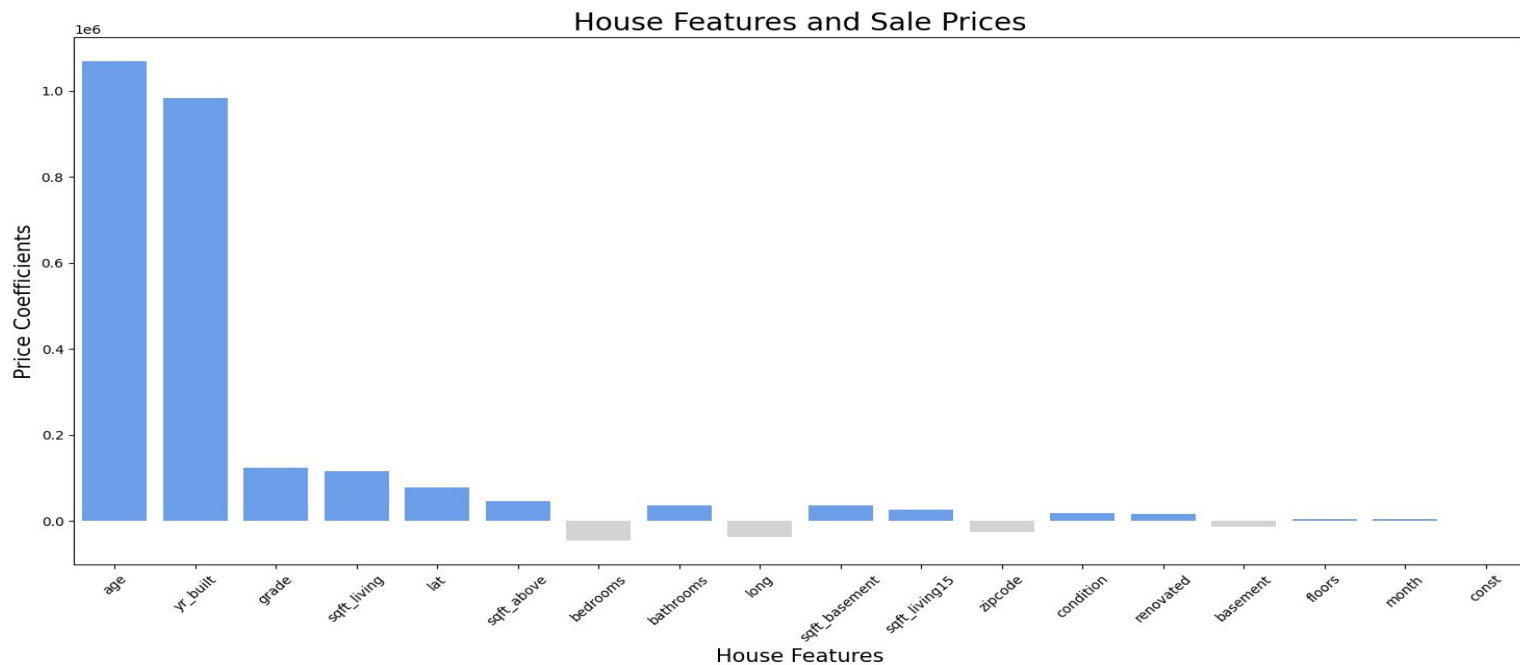
## Polynomial regression

The model has a good explanatory power.Involved transforming features into higher order polynomial terms to model nonlinear relationships.
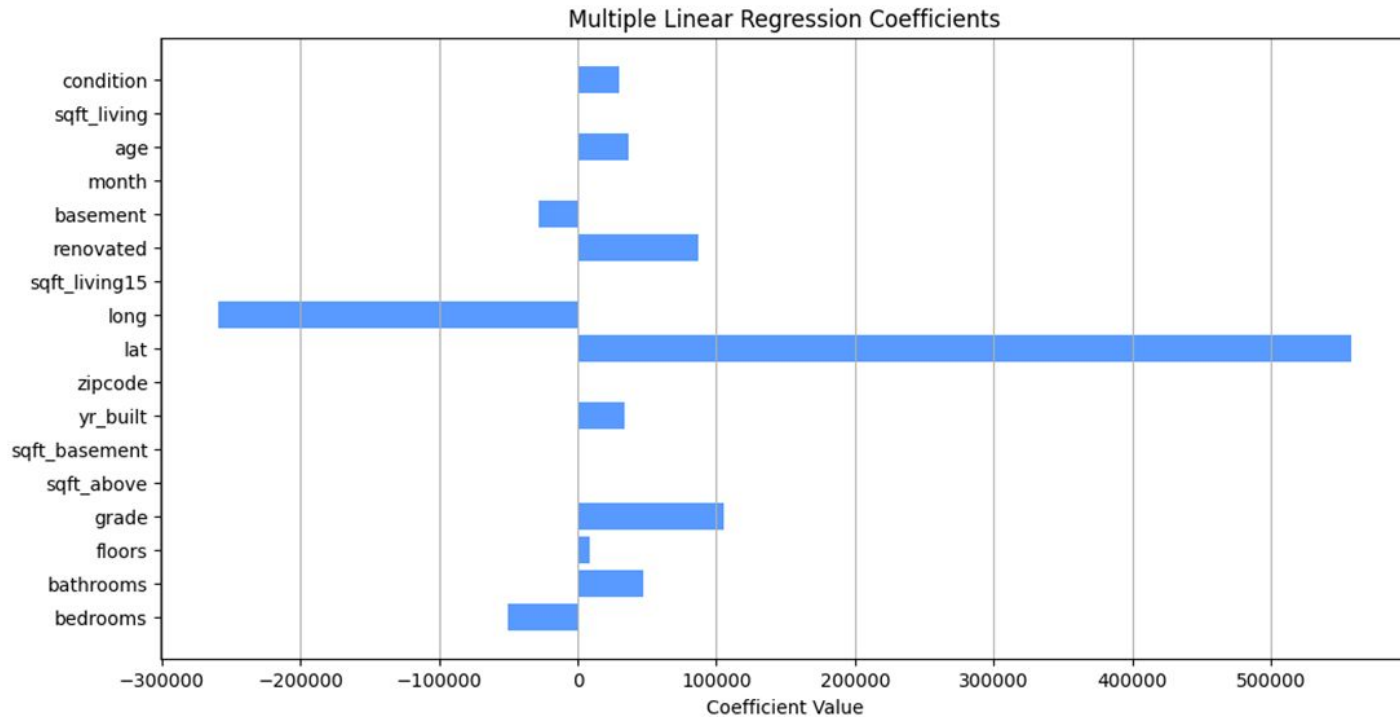
## Log Transformation

This model aimed to linearize relationships and stabilize the variance.It has improved model performance.
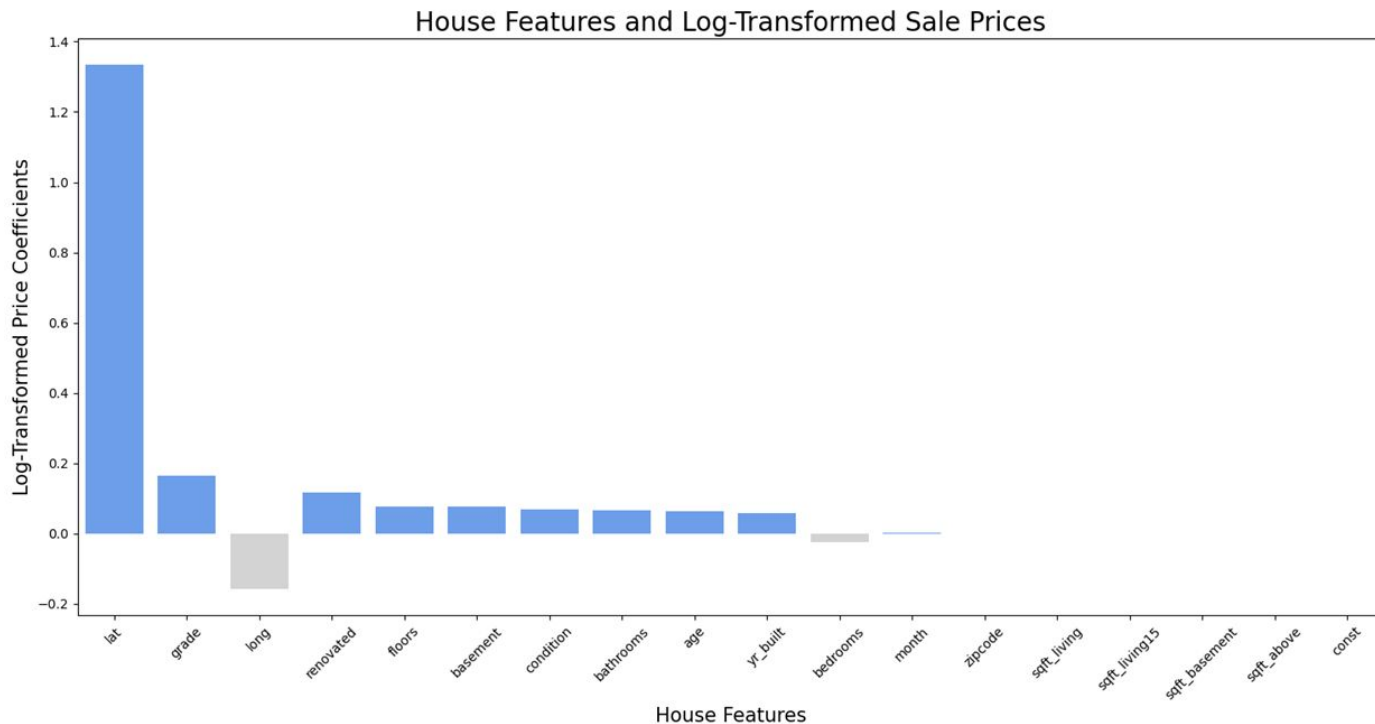
House Features and Sale Prices

The R-squared is approximately 0.668, suggesting that the model explains around 66.8% of the variance in home prices. Positive coefficients (e.g., sqft_living, bathrooms, grade) suggest an increase in these features corresponds to an increase in home price.
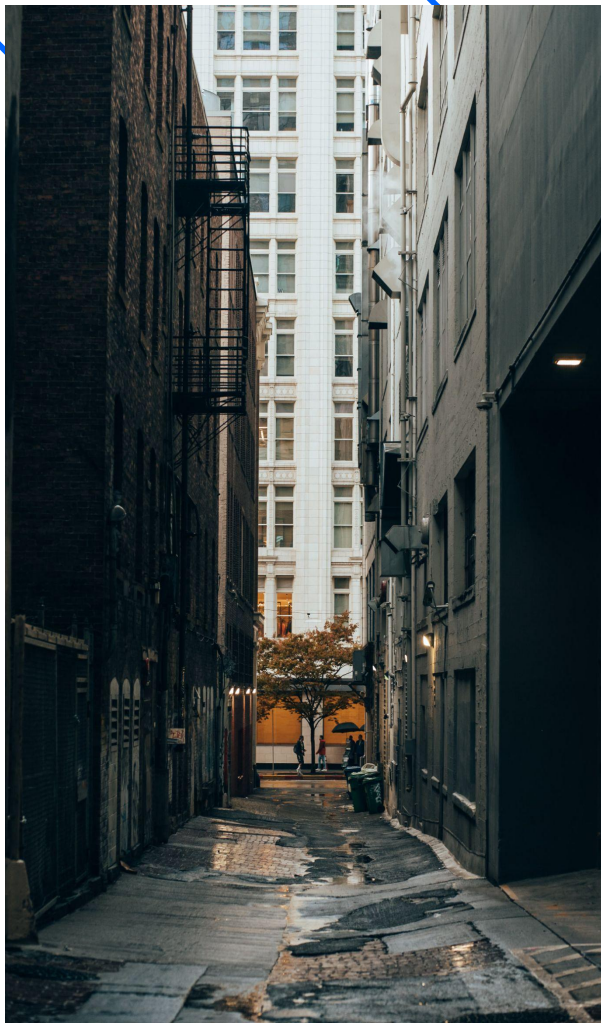
Multiple Linear Regression Coefficients

Bar graph illustrating the impact of features on the dependent variable in a multiple linear regression model, where positive and negative bar heights denote the direction and magnitude of influence, providing a concise overview of variable significance.

House Features and Log-Transformed Sale Prices

The R-squared is approximately 0.761, suggesting that the model explains around 76.1% of the variance in home prices. Positive coefficients (e.g., renovated, floors, grade) suggest an increase in these features corresponds to an increase in home price.
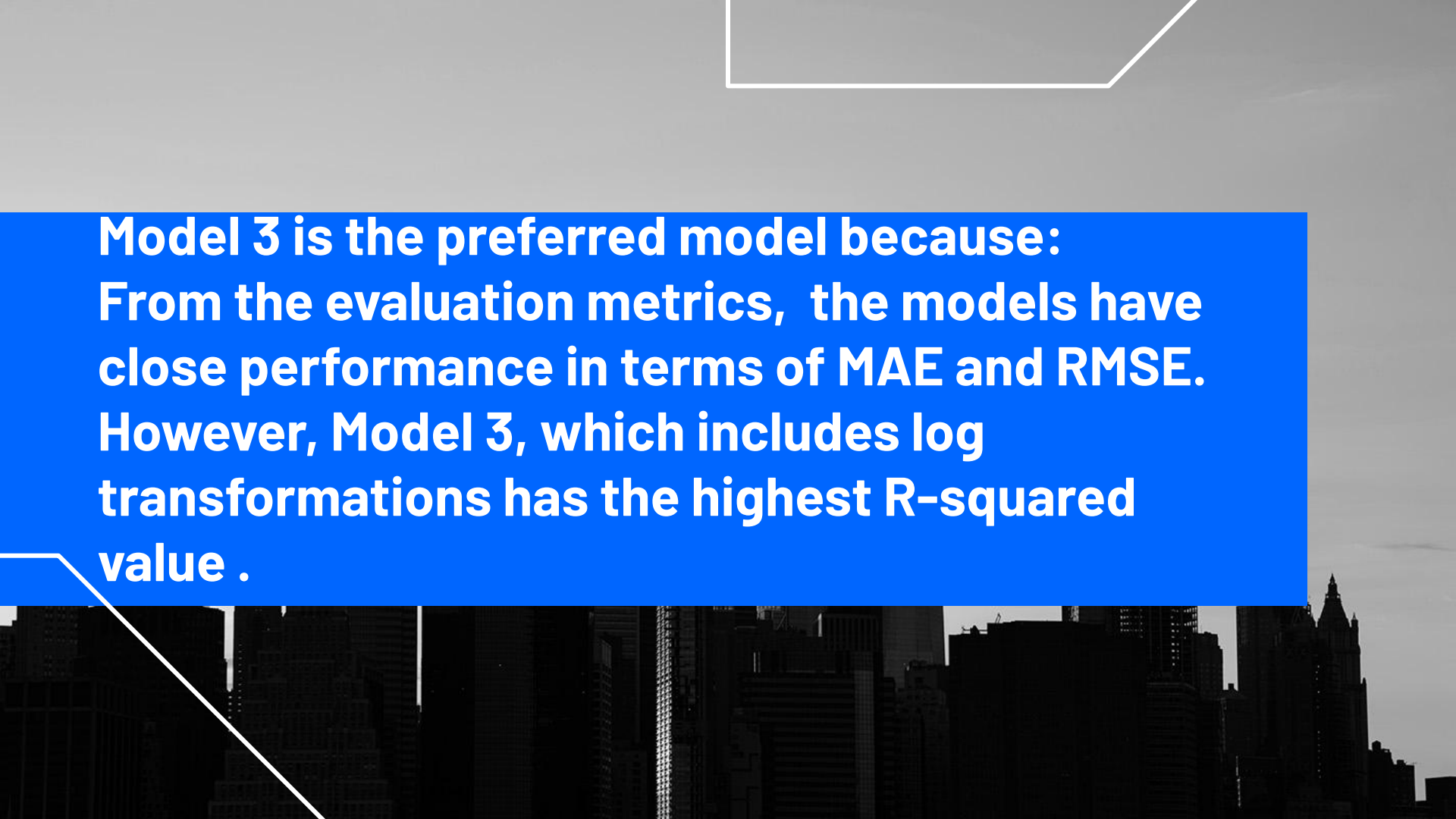
# Findings

- The log transformation of the variable has improved the models performance as indicated By the low RMSE and higher R squared .
- It has helped capture the underlying patterns in the data more effectively.

# MODEL COMPARISON

|  | MSE | R-Squared | F-stat |
|---|---|---|---|
| **Baseline** | 54.1 b | 0.668 | 4333 |
| **Polynomial** | 40.2 b | 0.663 | 1985 |
| **Log Trans** | 35.1 b | 0.761 | 4637 |

**Model 3 is the preferred model because:**
From the evaluation metrics, the models have close performance in terms of MAE and RMSE. However, Model 3, which includes log transformations has the highest R-squared value .

# KEY TAKEAWAYS

As the size of the living space increases, the estimated price of the house also increases. This indicates that larger houses are generally priced higher.

Older houses fetch higher prices. This could be due to factors such as historical significance or architectural value associated with older houses.

The more the bedrooms the more expensive the house, and he more space/land a house occupies,the more expensive it is

As the size of the living space increases, the estimated price of the house also increases. This indicates that larger houses are generally priced higher.

# CONCLUSION

# RECOMMENDATIONS

**Feature Enhancement**

Upgrade features like living space, property grade, or bathrooms to positively impact house prices

**Data Collection**

Collect additional data on location-specific factors, amenities, property age, and neighborhood characteristics for a more accurate regression model.

**Market Segmentation**

Analyze correlations between independent variables and house prices to identify market segments, such as luxury properties associated with higher grades
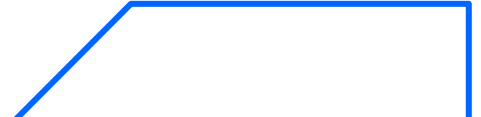
# LIMITATIONS

- **Non-Linearity**: Multiple linear regression struggles with non-linear relationships, leading to inaccurate predictions.
- **Causation vs. Correlation**: Correlation doesn't imply causation; the model identifies associations but can't establish cause-and-effect links.
- **Independence Assumption:** The model assumes independent observations; a violation can result in biased standard errors.
- **Linearity Assumption:** The model assumes a linear relationship; if not met, predictions may be biased, and coefficients inaccurate.

# NEXT STEPS

- **Alternative Models**: Explore complex models like random forests or neural networks to capture non-linear relationships effectively.
- **Time Trends and Seasonality:** If relevant, use time series models or include time-related features to better capture temporal patterns in the data.
- **Regular Monitoring:** Implement periodic model updates to incorporate new data, ensuring ongoing relevance and effectiveness.

# THANK YOU!

The Outliers