

# Fireplace Count Prediction Using KNN and MV Linear Regression

Ram Sunder

On 7th April 2025



# Table of Content

1. Overview
2. Understanding Dataset
3. Potential Issues to Address
4. Libraries used for EDA
5. Challenges(Future Warning)
6. Data Visualization
7. Summary
8. Question and Answers



# Understanding Dataset

## Initial Observations from the Dataset



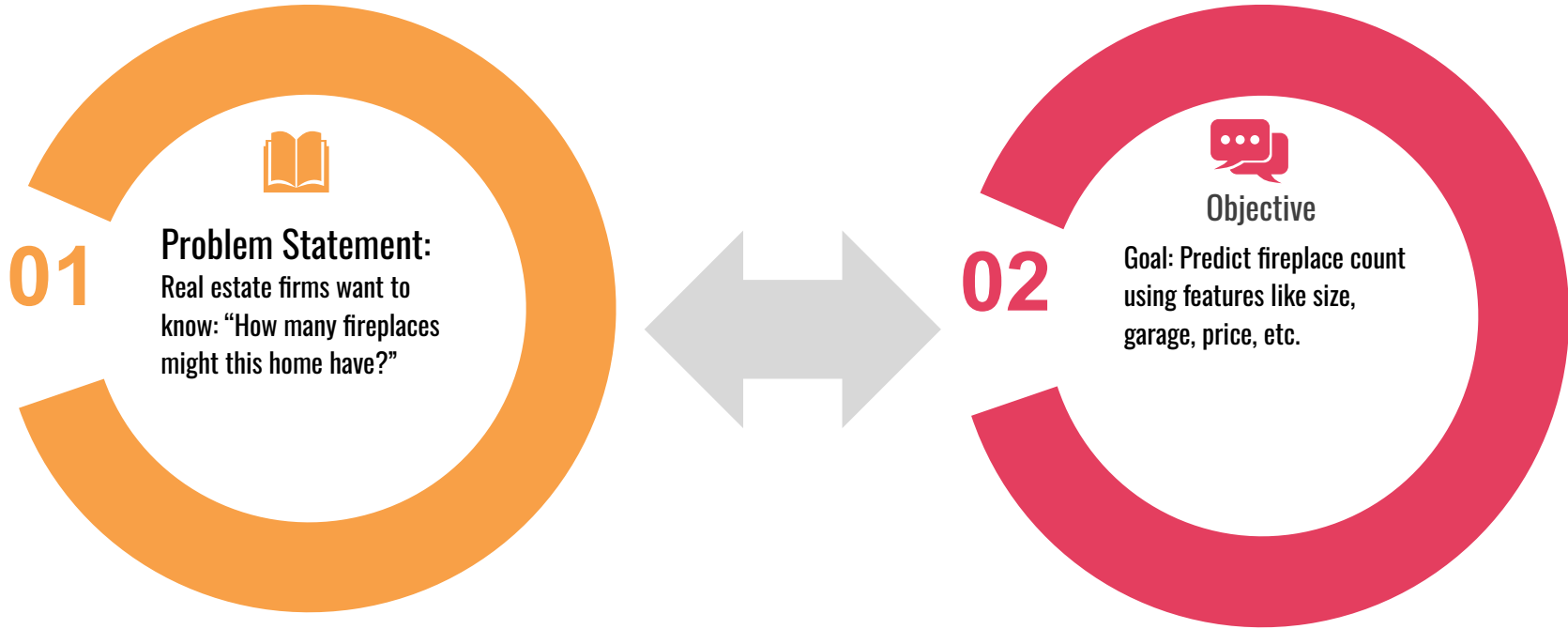
Total Records: 5,000 rows

Columns: 16

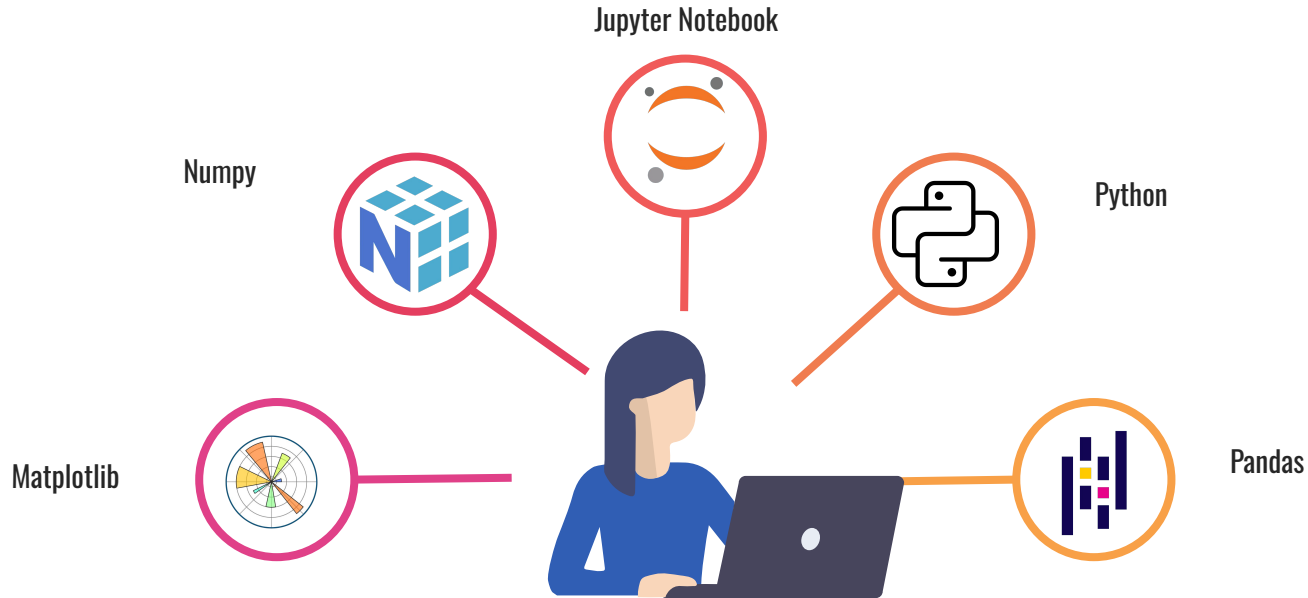
Data Types: Numerical: MLS, sold\_price, zipcode, longitude, latitude, lot\_acres, taxes, year\_built, bedrooms, fireplaces, bathrooms, garage, sqrt\_ft  
Data Types: Categorical (Objects): kitchen\_features, floor\_covering, HOA

Missing Values:  
lot\_acres (10 missing)  
fireplaces (25 missing)  
sqrt\_ft, garage, and HOA have None values (Need conversion to NaN)

# Potential Issues to Address



# Technology used in EDA



# Exploratory Data Analysis (EDA)

Overview of EDA: EDA is a crucial step in understanding the dataset and its underlying patterns before building any machine learning models.

## Missing Values

Identify missing values in the dataset.

## Summary Statistics

Analyze mean, median, and standard deviation for numerical variables.

## Data Distribution:

Visualize the distribution of key features (e.g., sold\_price, lot\_acres, etc.).



Outliers Handling: Mention any outliers that were identified and handled (if applicable).



Data Transformation: Explain any transformation, such as creating new features

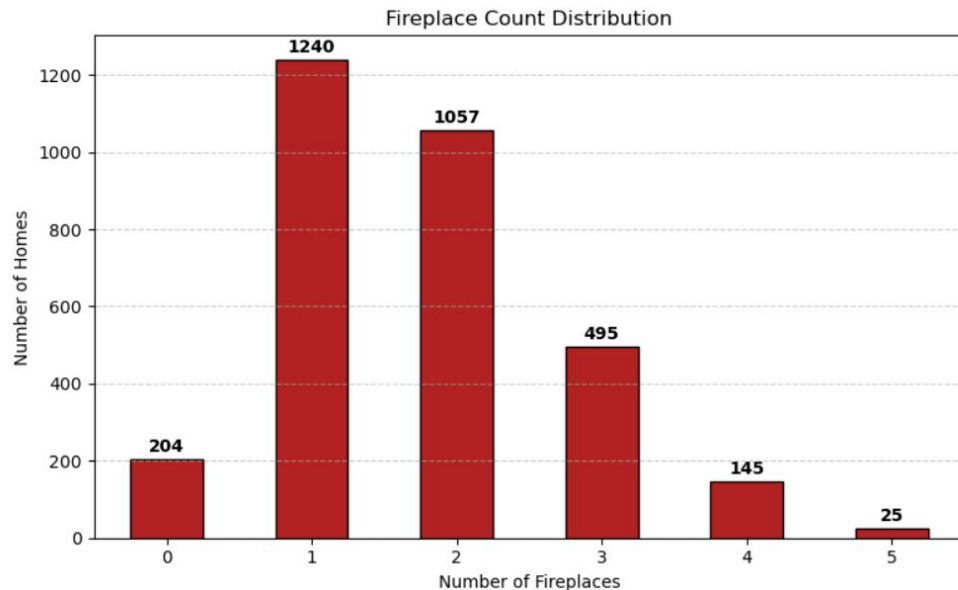


# Target Variable Creation:

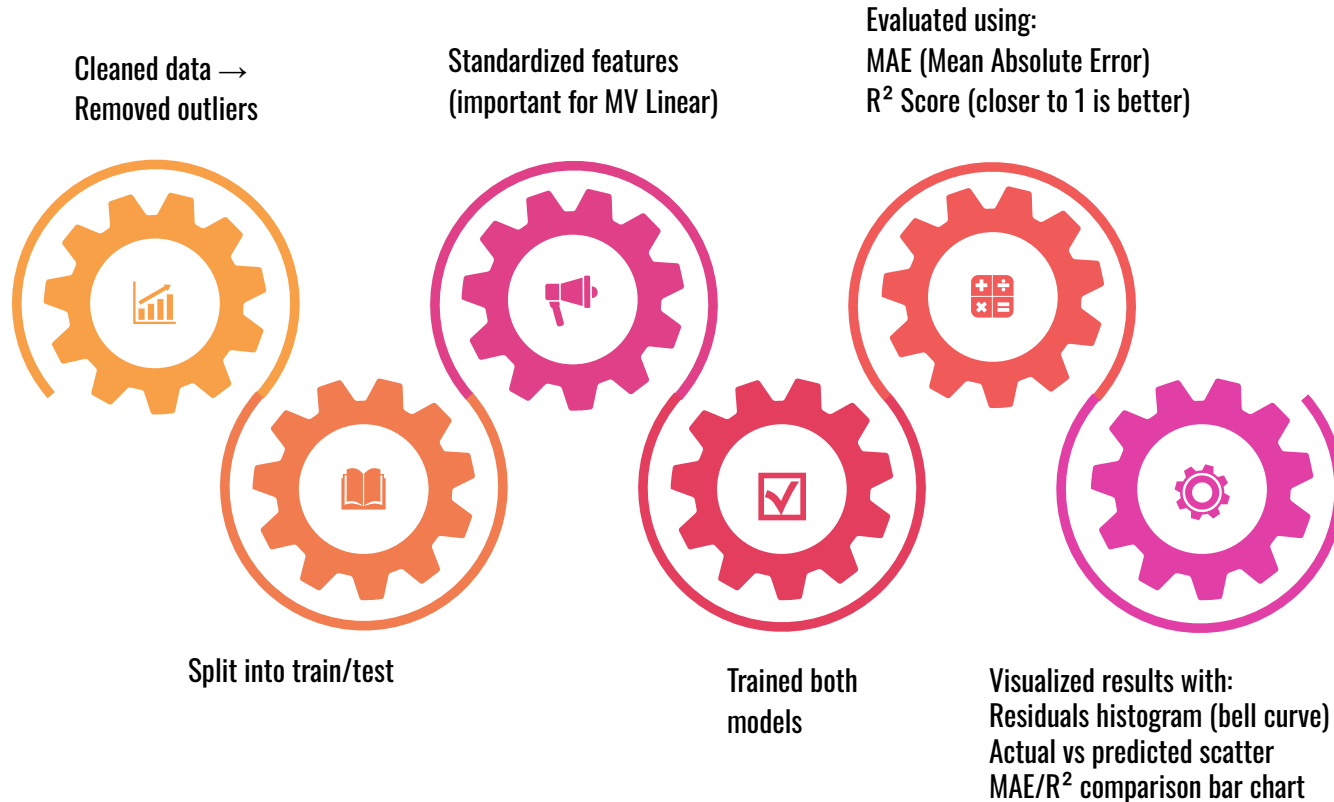
Target: **fireplaces** (0 to 5)

Features used:

- sold\_price, zipcode, sqrt\_ft, lot\_acres, year\_built, bedrooms, bathrooms, garage
- ~3166 cleaned row



# Implementation Flow





# Models Used

## A. KNN Regressor (Custom Implementation)

- Predicts based on similar nearby homes
- No assumption of linearity — great for real-world, uneven data



## B. Multivariate Linear Regression (Custom SGD version)

- Assumes a linear relationship between features and fireplaces
- We trained it using gradient descent

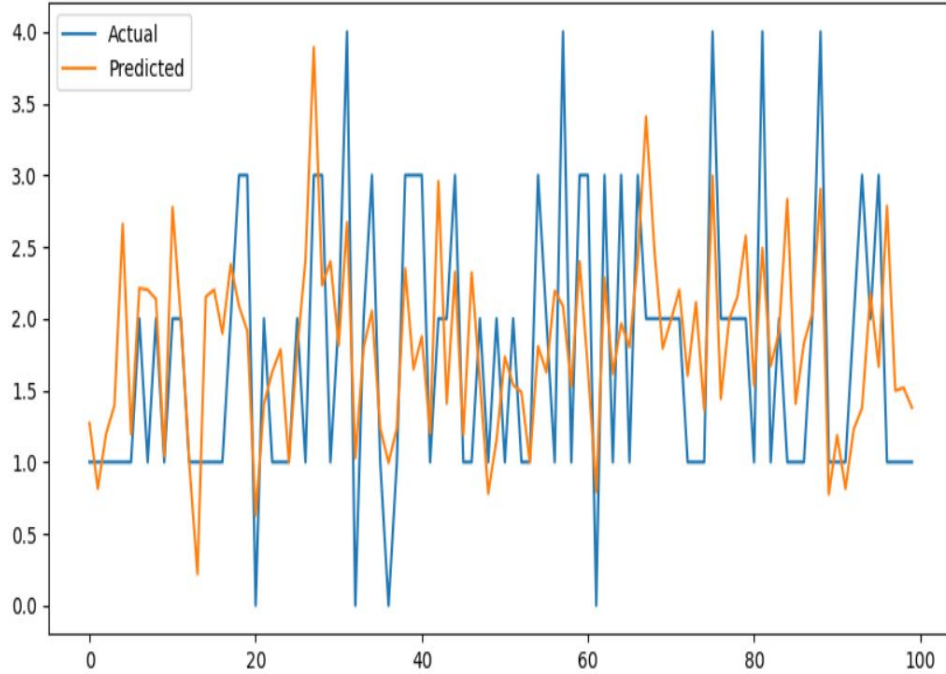
# Key Observations

- **KNN Regressor performed better**
  - Lower MAE
  - Positive  $R^2$  (closer to 0.5 depending on K)
- **Linear Regression failed**
  - $R^2$  was **negative**
  - Predictions were worse than just guessing the average
- Likely reason: **Fireplace counts are skewed** (mostly 1 or 2) → not ideal for linear models

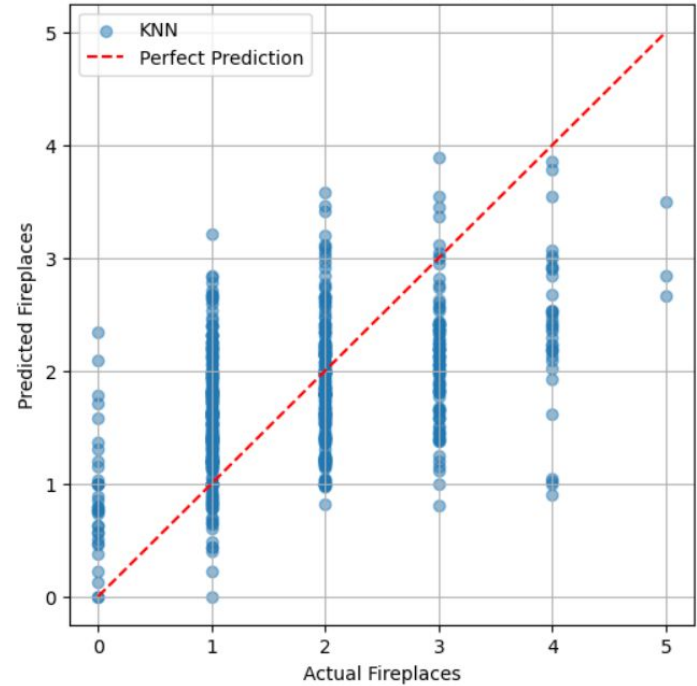


# KNN Regressor - Actual vs Predicted Fireplaces

KNN Regressor - Actual vs Predicted Fireplaces

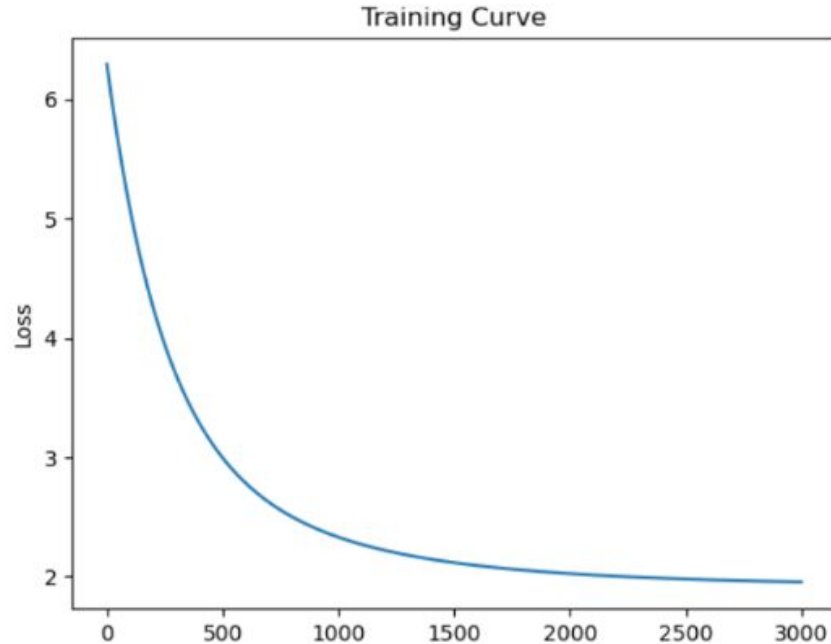


KNN: Actual vs Predicted

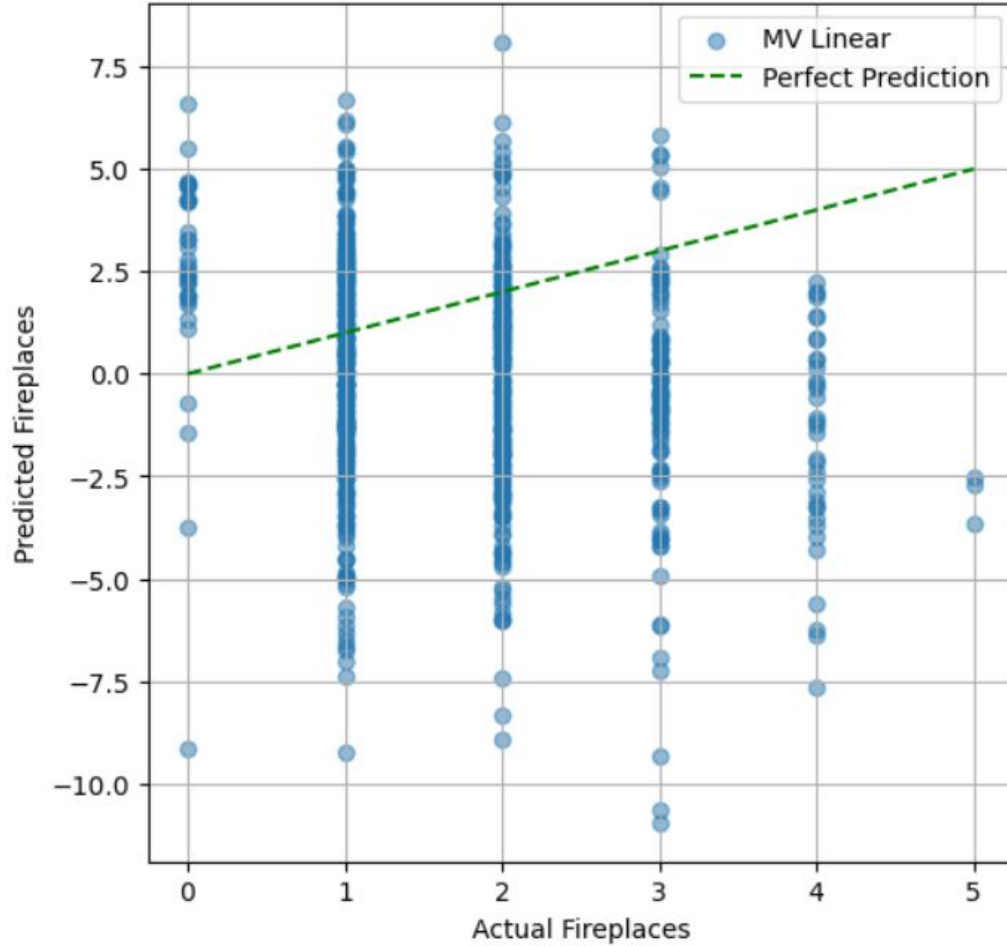


# Train and Predict with MV Linear Regression

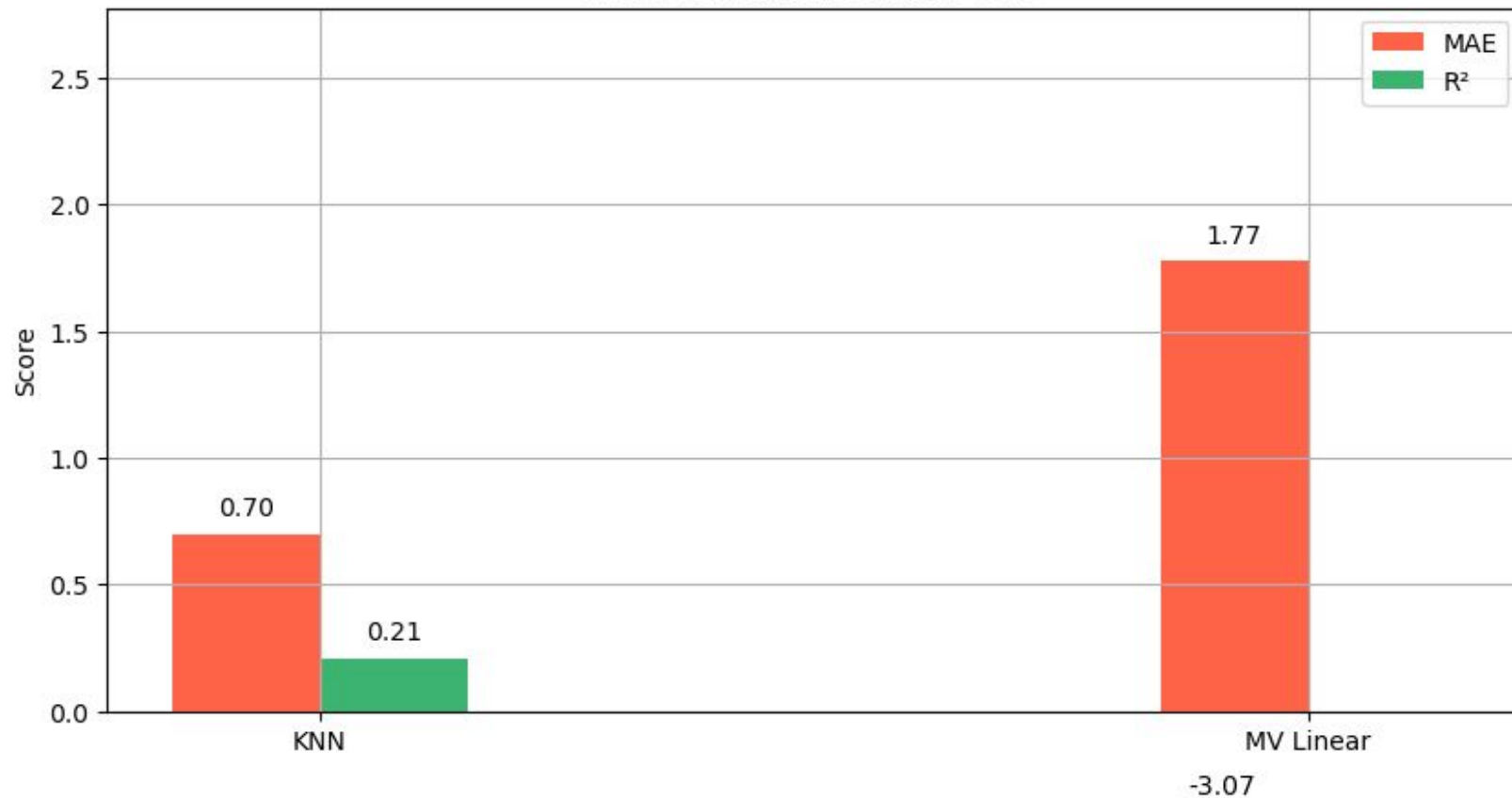
```
linreg = MVLinearRegression()  
linreg.fit(X_train_scaled, y_train, eta=1e-3, epochs=3000)  
y_pred_lin = linreg.predict(X_test_scaled)
```



MV Linear: Actual vs Predicted



Model Performance: MAE vs  $R^2$



# Improvements / Future Work

- Try KNN Classifier or Decision Trees (fireplaces = discrete)
- Add more features (e.g., neighborhood quality, presence of chimney)
- Use GridSearch for better K in KNN
- Handle skew with techniques like SMOTE (for classification)
- Try ensemble models like Random Forests for more power

# Final Statement

- Built and compared two custom ML models to predict the number of fireplaces in homes.
- KNN performed better due to the non-linear nature of the data, while MV Linear struggled.



**Thank you!**