# End-to-End Machine Learning Pipeline for Real Estate Valuation

**Nawaraj Paudel, PhD**
**(Quantitative Modeling of Materials)**
**Data Scientist and ML Engineer**

# Real estate tech: Where massive markets meet innovation

## Real Estate Market by the Numbers

➢ 146 million residential units valued at 43 trillion USD

➢ Commercial real estate valued at 21 trillion USD

➢ 2- 8 % of these properties sold every year

## Marketcap Comparison

➢ S &P 500 : 45 Trillion

➢ NASDAQ 100 : 20 Trillion

## Tech Transformation: Key Real Estate Domains

➢ PropTech: Real estate markets

➢ ConTech: Construction startups

➢ SmartRealEstate: Intelligent cities and buildings

➢ RealEstateFinTech: Mortage marketplace, Blockchain and smart contracts, Crowdfunding platforms

➢ Collaborative economy

# The housing data contains 80 features including 43 categorical features

## Dwelling Characteristics

➢ **OverallQual:** Rates the overall material and finish of the house

➢ **YearBuilt:** Original construction date

➢ **Year Remod/Add:** Remodel date

## Living Area

➢ **1st Flr SF**: First Floor square feet

➢ **2nd Flr SF**: Second floor square feet

➢ **Low Qual Fin SF**: Low quality finished square feet (all floors)

➢ **Gr Liv Area**: Above grade (ground) living area square feet

## Bedrooms and Bathrooms

➢ **Bsmt Full Bath:** Basement full bathrooms

➢ **Bsmt Half Bath:** Basement half bathrooms

➢ **Full Bath:** Full bathrooms above grade

➢ **Half Bath:** Half baths above grade
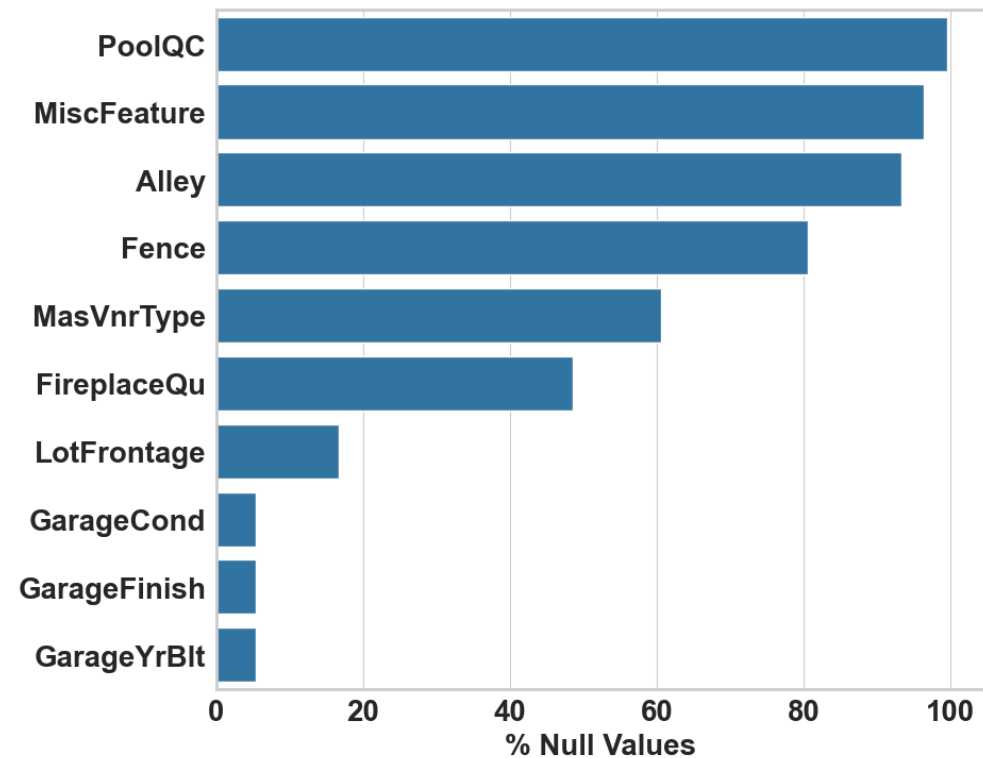
➢ **Bedroom:** Bedrooms above grade

## Other Features

➢ **Fireplaces:** Number of fireplaces

➢ **FireplaceQu:** Fireplace quality

➢ **Garage Cars:** Size of garage in car capacity

➢ **Garage Area:** Size of garage in square feet

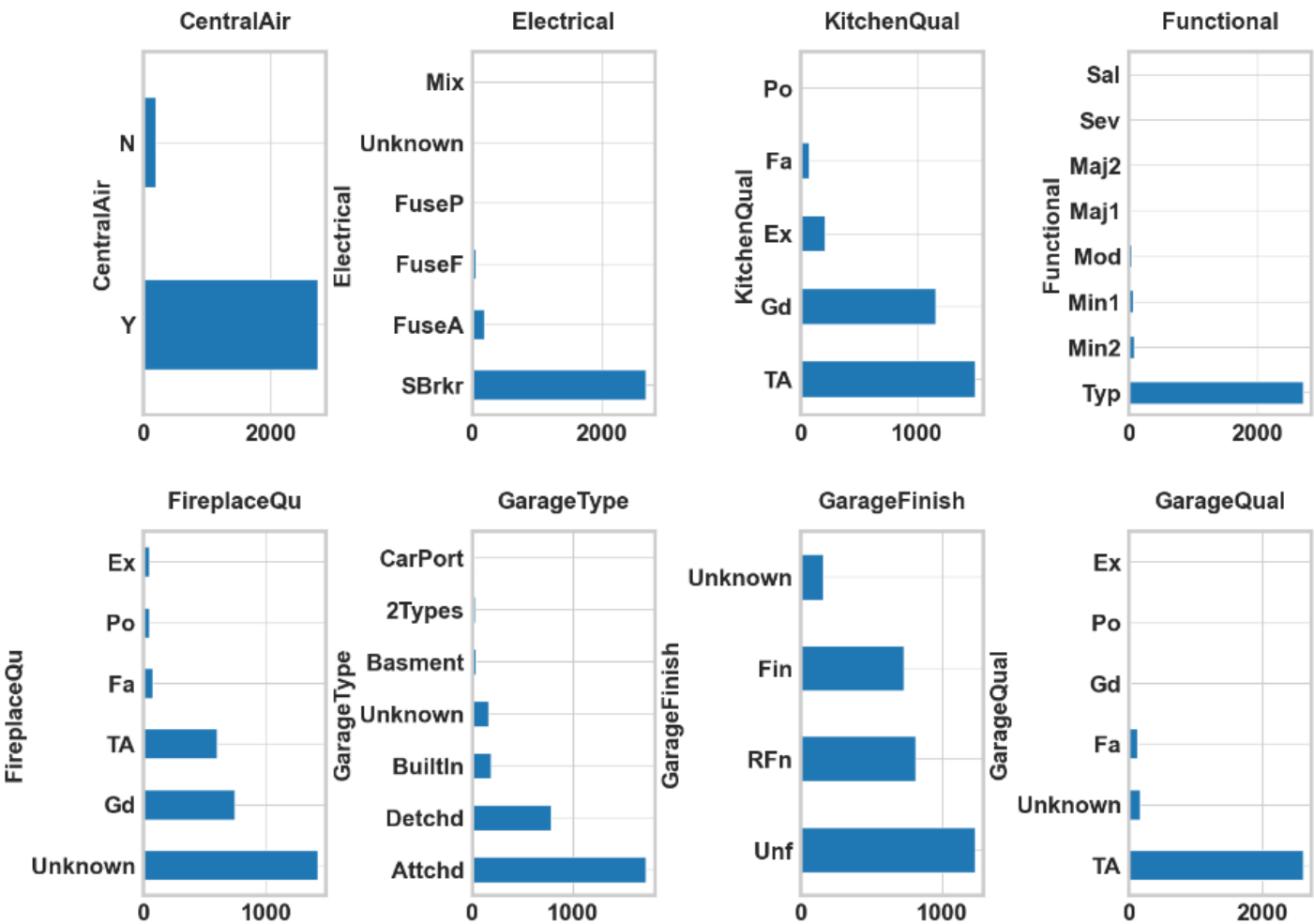➢ 64 other features

Property Recommendation

Sale Price Prediction

# Missing data were labeled 'unknown' for categorical features

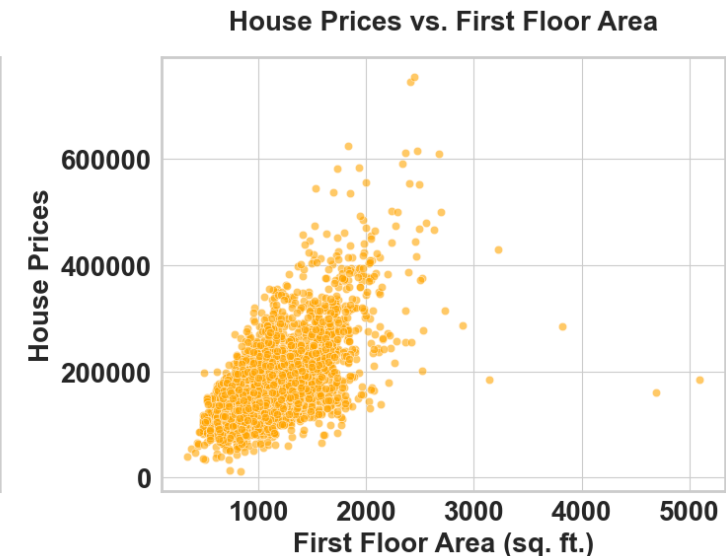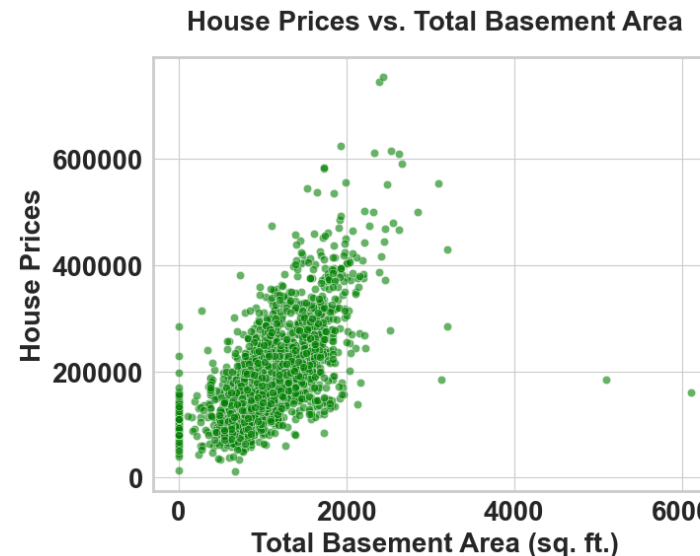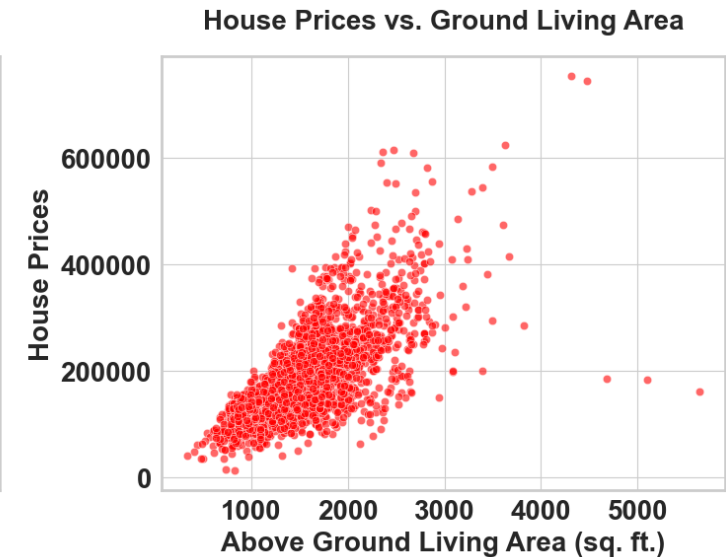**Top 10 Features with the Highest Percentage of Null Values**



If a categorical column had more than 10% missing values, a new category called 'Unknown' was created. For columns with less than 10% missing values, they were imputed with the most frequent category.

# House prices have linear relationship with area and quality

- ➢ The house of price increases with quality and area of the house

- ➢ The rate of change of house price is steeper for total basement area

- ➢ There are price ranges for same quality and area stemming a fan like structure more pronounced in 'AboveGround Living Area' and 'First Floor Area'

- ➢ There are some outliers in area features

# Q2/Q3: Peak season for home sales

## Number of Houses Sold by Month and Year



**House Sales Seasonal Trend**

➢ Sales start increasing in May and peak in June

➢ Lowest transactions occur at the end of the year and extend into the first two months of the following year

# Neighborhood-based median house prices



Sales Price of Homes Available by Top 10 Neighborhood in Ames City

- House prices in some neighborhoods start from as low as $200K, while the highest-priced neighborhoods can reach up to $600K

- Median house prices vary depending on the neighborhood

- Certain neighborhoods exhibit outlier house prices, with some properties significantly deviating from the median range

# Engineered features are highly correlated with original features

**Engineered Features**
- TotalBaths
- HouseAge
- YearRemodAge
- TotalSqFt

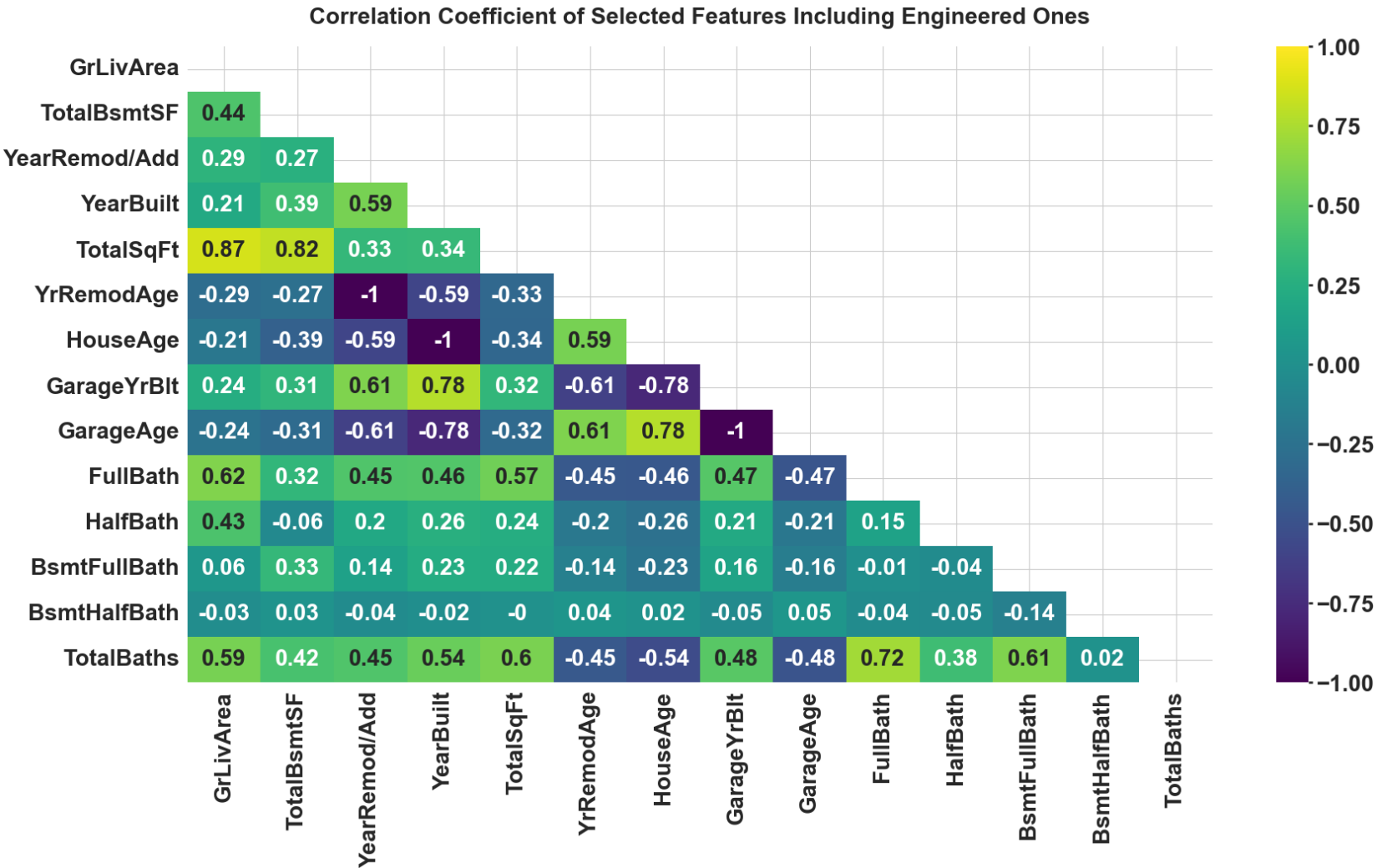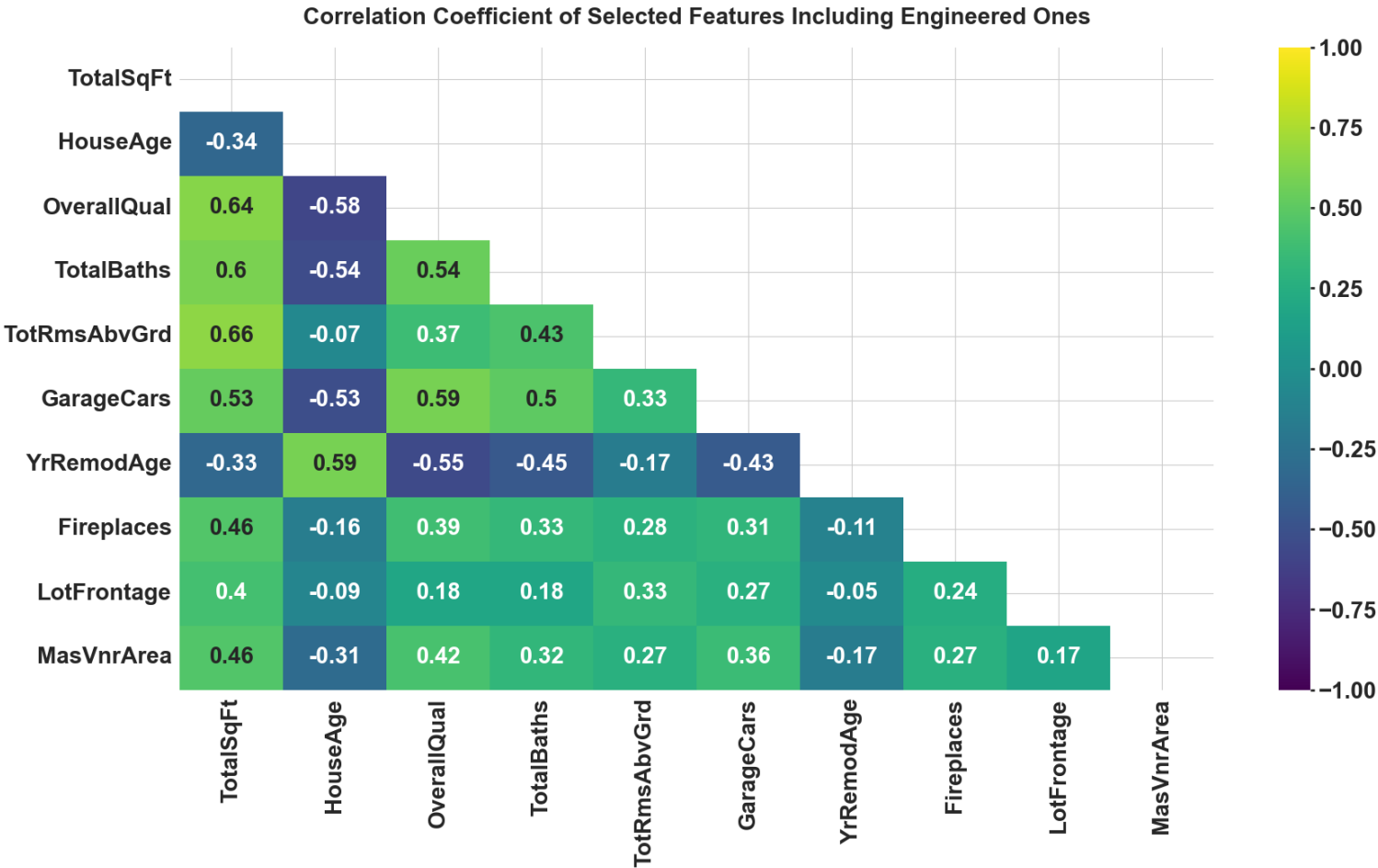| Features | VIF |
|---|---|
| TotalSqFt | 3201.15 |
| BsmtFinSF1 | 1021.55 |
| BsmtUnfSF | 954.11 |
| 2ndFlrSF | 915.30 |
| 1stFlrSF | 755.71 |
| BsmtFinSF2 | 145.64 |
| LowQualFinSF | 11.83 |
| GarageCars | 5.60 |
| GarageArea | 5.44 |
| TotRmsAbvGrd | 4.50 |
| HouseAge | 4.49 |
| GarageAge | 3.20 |
| TotalBaths | 3.04 |
| YrRemodAge | 2.37 |

**Correlation Coefficient of Selected Features Including Engineered Ones**

| | GrLivArea | TotalBsmtSF | YearRemod/Add | YearBuilt | TotalSqFt | YrRemodAge | HouseAge | GarageYrBlt | GarageAge | FullBath | HalfBath | BsmtFullBath | BsmtHalfBath | TotalBaths |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GrLivArea | | | | | | | | | | | | | | |
| TotalBsmtSF | 0.44 | | | | | | | | | | | | | |
| YearRemod/Add | 0.29 | 0.27 | | | | | | | | | | | | |
| YearBuilt | 0.21 | 0.39 | 0.59 | | | | | | | | | | | |
| TotalSqFt | 0.87 | 0.82 | 0.33 | 0.34 | | | | | | | | | | |
| YrRemodAge | -0.29 | -0.27 | -1 | -0.59 | -0.33 | | | | | | | | | |
| HouseAge | -0.21 | -0.39 | -0.59 | -1 | -0.34 | 0.59 | | | | | | | | |
| GarageYrBlt | 0.24 | 0.31 | 0.61 | 0.78 | 0.32 | -0.61 | -0.78 | | | | | | | |
| GarageAge | -0.24 | -0.31 | -0.61 | -0.78 | -0.32 | 0.61 | 0.78 | -1 | | | | | | |
| FullBath | 0.62 | 0.32 | 0.45 | 0.46 | 0.57 | -0.45 | -0.46 | 0.47 | -0.47 | | | | | |
| HalfBath | 0.43 | -0.06 | 0.2 | 0.26 | 0.24 | -0.2 | -0.26 | 0.21 | -0.21 | 0.15 | | | | |
| BsmtFullBath | 0.06 | 0.33 | 0.14 | 0.23 | 0.22 | -0.14 | -0.23 | 0.16 | -0.16 | -0.01 | -0.04 | | | |
| BsmtHalfBath | -0.03 | 0.03 | -0.04 | -0.02 | -0 | 0.04 | 0.02 | -0.05 | 0.05 | -0.04 | -0.05 | -0.14 | | |
| TotalBaths | 0.59 | 0.42 | 0.45 | 0.54 | 0.6 | -0.45 | -0.54 | 0.48 | -0.48 | 0.72 | 0.38 | 0.61 | 0.02 | |

# The top 10 numerical features were selected for modeling, ensuring the most impactful variables are used for accurate predictions

| Feature | F-Score | P-Value |
|---|---|---|
| OverallQual | 4342.9792 | 0.000000 |
| TotalSqFt | 3873.3856 | 0.000000 |
| GarageCars | 1752.9074 | 0.000000 |
| TotalBaths | 1709.5729 | 0.000000 |
| GarageArea | 1665.5998 | 0.000000 |
| 1stFlrSF | 1519.5015 | 0.000000 |
| HouseAge | 1053.7745 | 0.000000 |
| YrRemodAge | 939.3147 | 0.000000 |
| GarageAge | 817.7759 | 0.000000 |
| MasVnrArea | 806.9772 | 0.000000 |
| TotRmsAbvGrd | 736.0766 | 0.000000 |
| Fireplaces | 718.7684 | 0.000000 |
| BsmtFinSF1 | 543.6743 | 0.000000 |
| WoodDeckSF | 317.1991 | 0.000000 |
| LotFrontage | 280.3683 | 0.000000 |

| Variable | VIF |
|---|---|
| TotalSqFt | 3.55 |
| HouseAge | 3.42 |
| OverallQual | 2.77 |
| TotalBaths | 2.10 |
| YrRemodAge | 1.94 |
| TotRmsAbvGrd | 1.93 |
| GarageCars | 1.90 |
| Fireplaces | 1.41 |
| MasVnrArea | 1.36 |
| LotFrontage | 1.36 |



Correlation Coefficient of Selected Features Including Engineered Ones

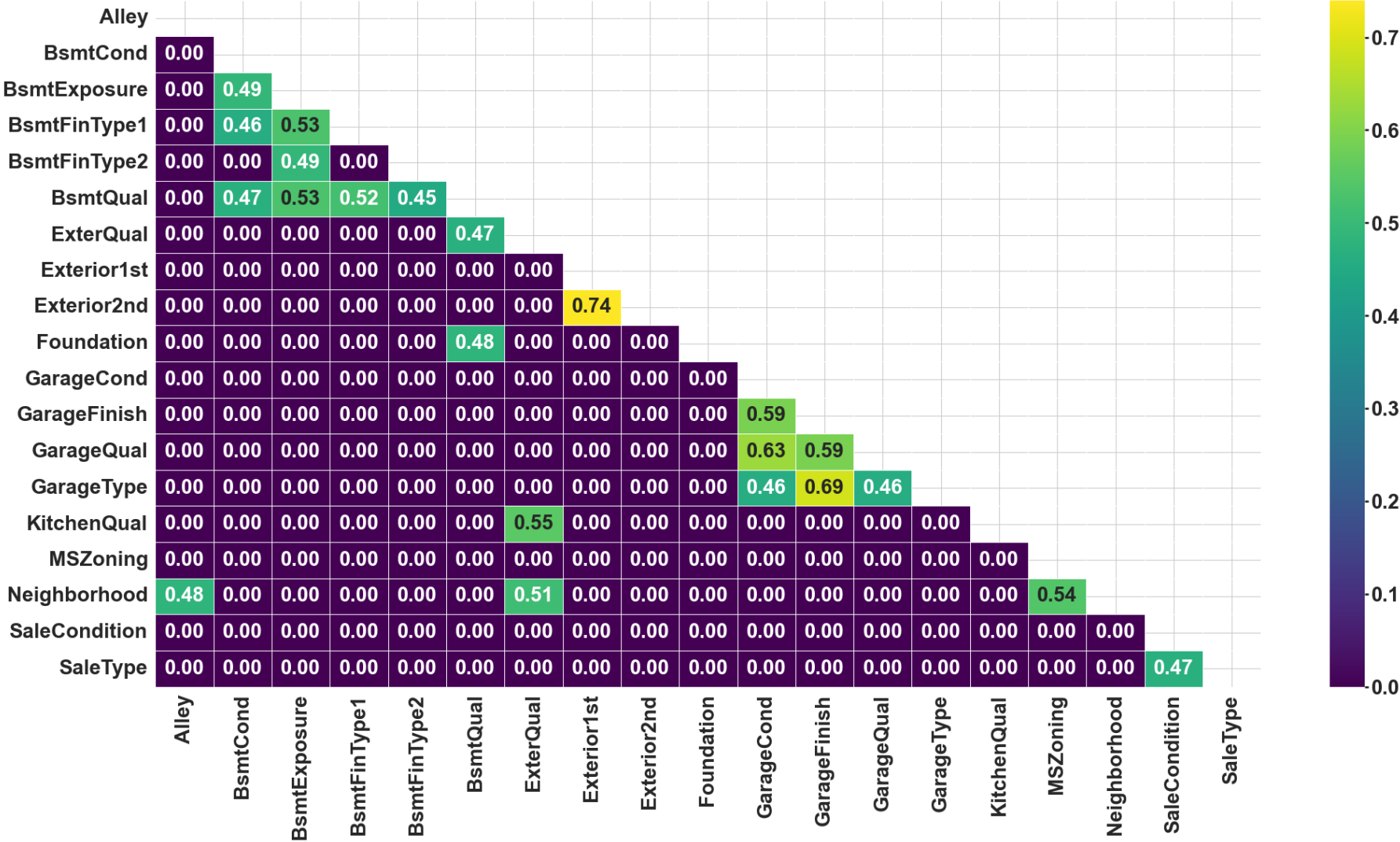# The top 4 categorical features were selected using ANOVA and Cramer's V

## ANOVA Test

| Categorical Column | F-value | P-value |
|---|---|---|
| ExterQual | 808.16 | 0 |
| KitchenQual | 555.20 | 0 |
| BsmtQual | 489.49 | 0 |
| GarageFinish | 341.35 | 1.30E-185 |
| FireplaceQu | 220.99 | 1.81E-195 |
| CentralAir | 194.12 | 1.52E-42 |
| Foundation | 186.77 | 3.43E-169 |
| HeatingQC | 161.99 | 2.40E-123 |
| MasVnrType | 133.61 | 1.70E-103 |
| GarageType | 130.70 | 4.59E-144 |
| BsmtExposure | 130.08 | 5.58E-101 |
| BsmtFinType1 | 117.58 | 5.55E-131 |
| Neighborhood | 117.23 | 0 |



Cramér's V Heatmap for Association between Categorical Features with V-value above 0.45

# Minimum viable product (MVP): CatBoost trained with all features

➤ TotalSqFt and OverallQual emerged as the most significant features, contributing over 40% to the model's predictive power

➤ Over 10 features individually contribute 1 - 2% to the model's predictions

➤ The CatBoost model, trained with default parameters and using all input features, achieved an average out-of-fold (OOF) R2 score of 9.125%
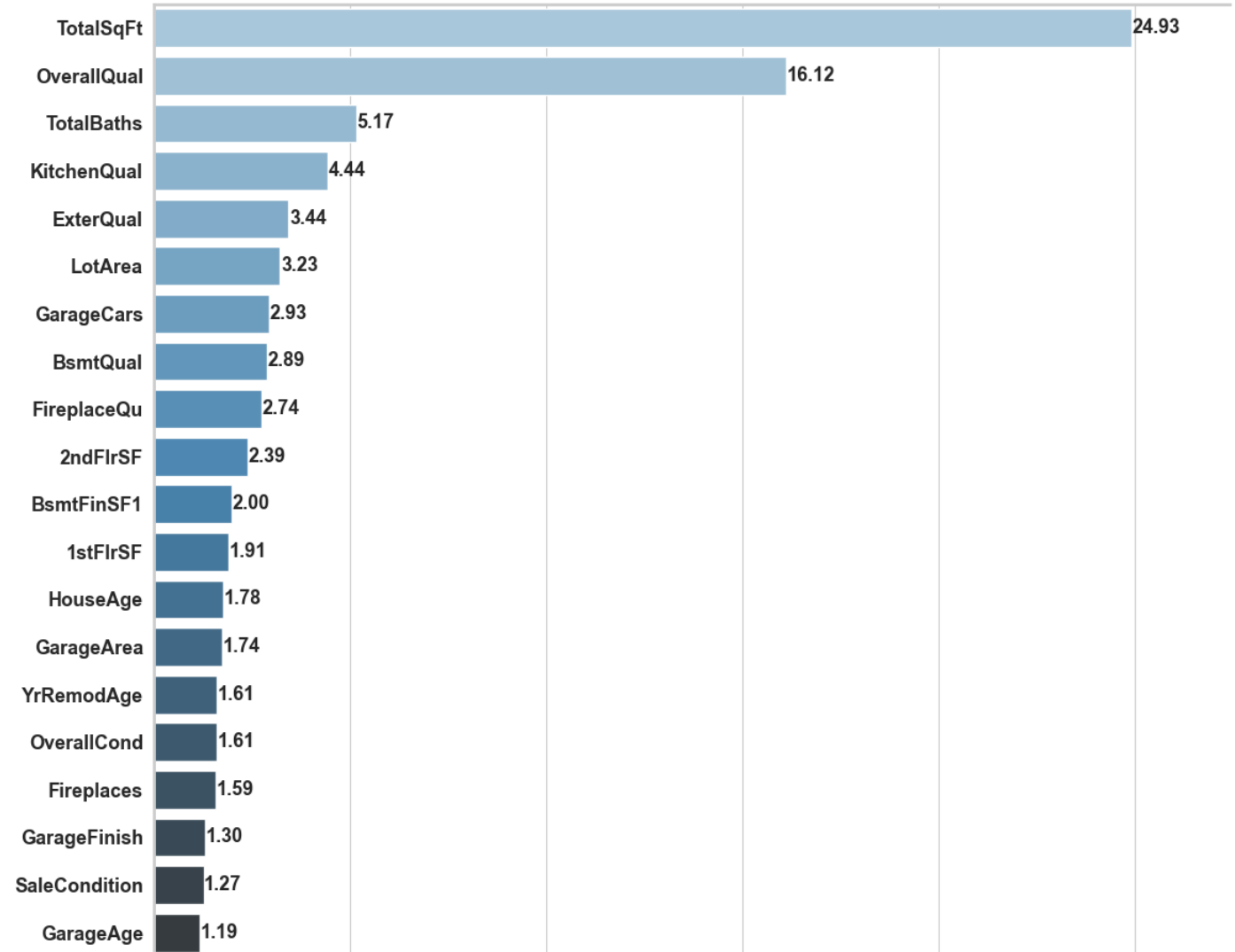
```
1  # Identify and fill NaNs in categorical columns
2  cat_features = [col for col in X_train_new.columns if X_train_new[col].dtype == 'object']
3  # Define and train the default CatBoost Model
4  base_model = CatBoostRegressor(cat_features = cat_features, random_state = 42, verbose = 0)
5  base_scores = cross_val_score(base_model, X_train_new, y_train, cv = 5, scoring = 'r2')
6  print(f"Average r2 score for default CatBoost: {base_scores.mean():.4f}")
[131]
...  Average r2 score for default CatBoost: 0.9125

1  base_scores
[132]
...  array([0.88449204, 0.92523497, 0.92133909, 0.89634955, 0.93528412])
```
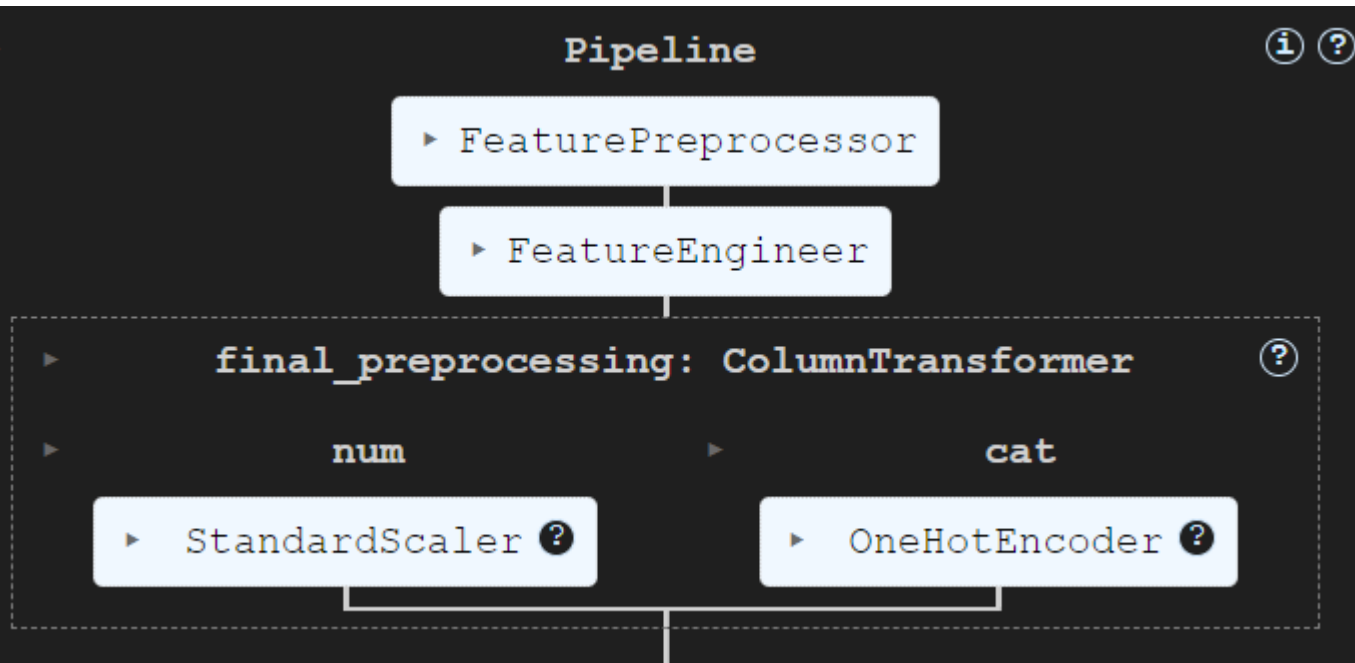
### Top 20 Most Important Features with Importance Scores - CatBoost Model

| Feature | Importance Score |
|---|---|
| TotalSqFt | 24.93 |
| OverallQual | 16.12 |
| TotalBaths | 5.17 |
| KitchenQual | 4.44 |
| ExterQual | 3.44 |
| LotArea | 3.23 |
| GarageCars | 2.93 |
| BsmtQual | 2.89 |
| FireplaceQu | 2.74 |
| 2ndFlrSF | 2.39 |
| BsmtFinSF1 | 2.00 |
| 1stFlrSF | 1.91 |
| HouseAge | 1.78 |
| GarageArea | 1.74 |
| YrRemodAge | 1.61 |
| OverallCond | 1.61 |
| Fireplaces | 1.59 |
| GarageFinish | 1.30 |
| SaleCondition | 1.27 |
| GarageAge | 1.19 |

# Automated preprocessing pipeline



**Feature Preprocessor**

➢ Data loading, column names normalization and data validation

➢ Handles missing values using median imputation for numeric features and 'Unknown' for categorical features

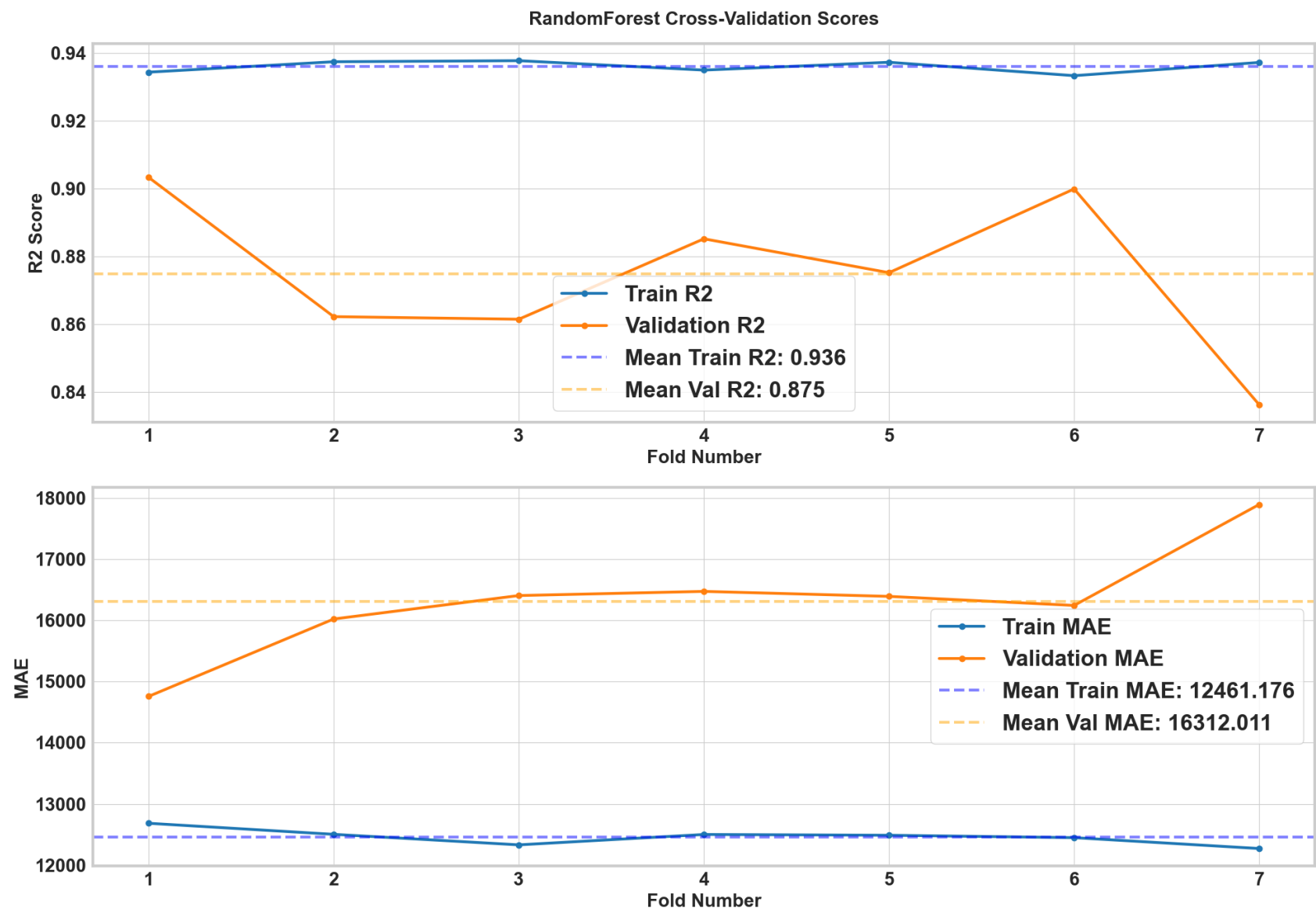➢ Consolidates rare categories based on threshold (merges categories < 8% into 'Other')

**Feature Engineer**

➢ Creates engineered features: TotalSqFt, HouseAge , TotalBaths, and YrRemodAge

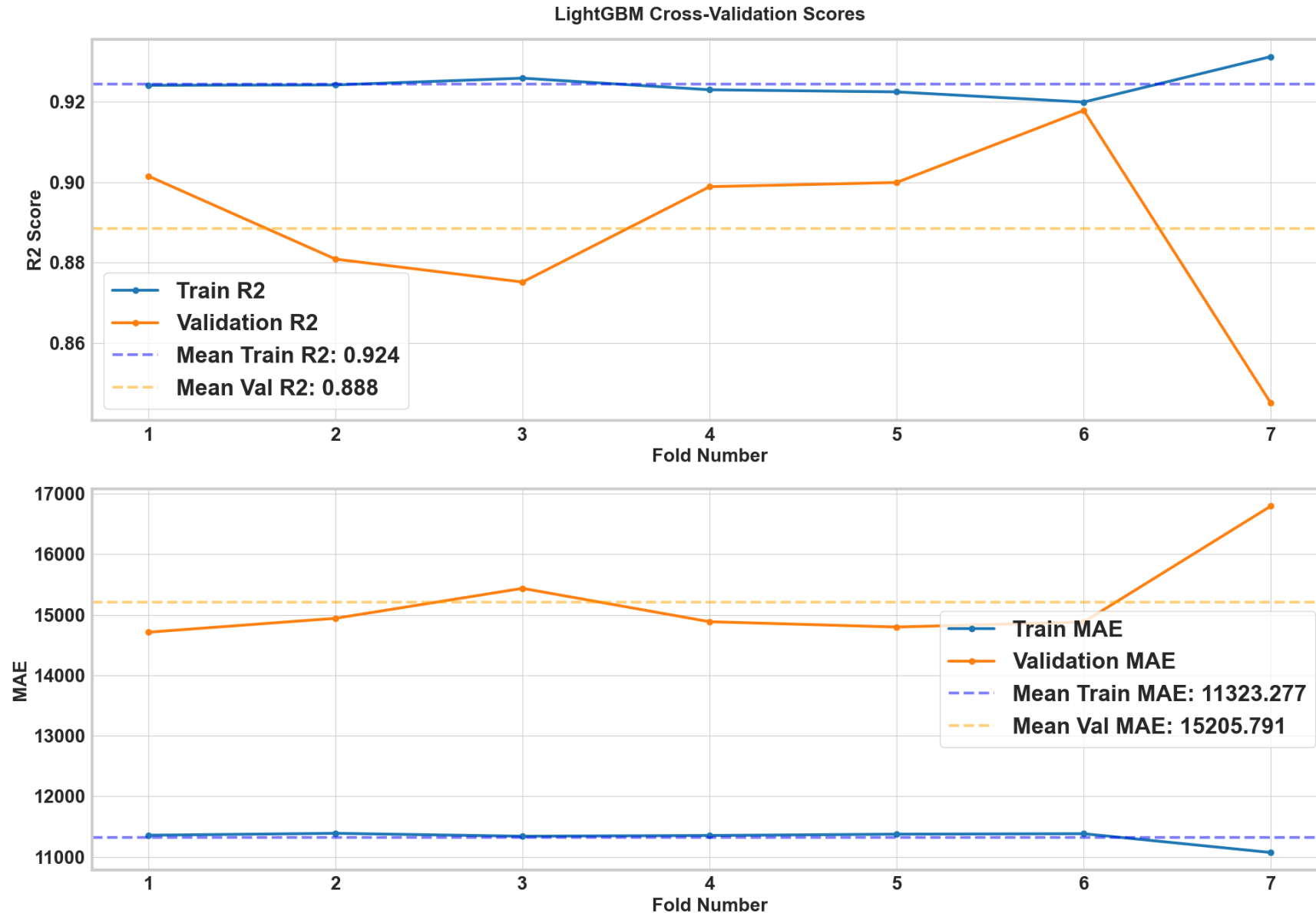➢ Automatically drops original features after engineering new ones

**Final Column Transformers**

➢ Numeric Pipeline: Applies StandardScaler to normalize all numeric features

➢ Categorical Pipeline: Uses OneHotEncoder with drop = 'first' and handle_unknown = 'ignore' for categorical variables

```
Pipeline(steps=[('initial_preprocessing',          FeaturePreprocessor(categorical_features=['Neighborhood',
'FireplaceQu',                                  'KitchenQual',
'BsmtExposure'],                       numeric_features=['YearRemodAdd', 'YrSold',
'Fireplaces',                          'LotFrontage',                          'GarageCars',
'MasVnrArea',                           'BsmtFullBath',                        'GrLivArea',
'BsmtHalfBath',                          'YearBuilt',                          'OverallQual',
'TotalBsmtSF',                          'HalfBath', 'FullBath',                        'TotRmsAb...
Pipeline(steps=[('scaler',                                StandardScaler())]),
['OverallQual',                        'TotRmsAbvGrd', 'GarageCars',
'Fireplaces', 'LotFrontage',                          'MasVnrArea', 'TotalSqFt',
'HouseAge', 'TotalBaths',                          'YrRemodAge']),                          ('cat',
Pipeline(steps=[('onehot',                                OneHotEncoder(drop='first',
handle_unknown='ignore',                                sparse_output=False)]),
['Neighborhood',                        'FireplaceQu', 'KitchenQual',
'BsmtExposure'])])])])
```
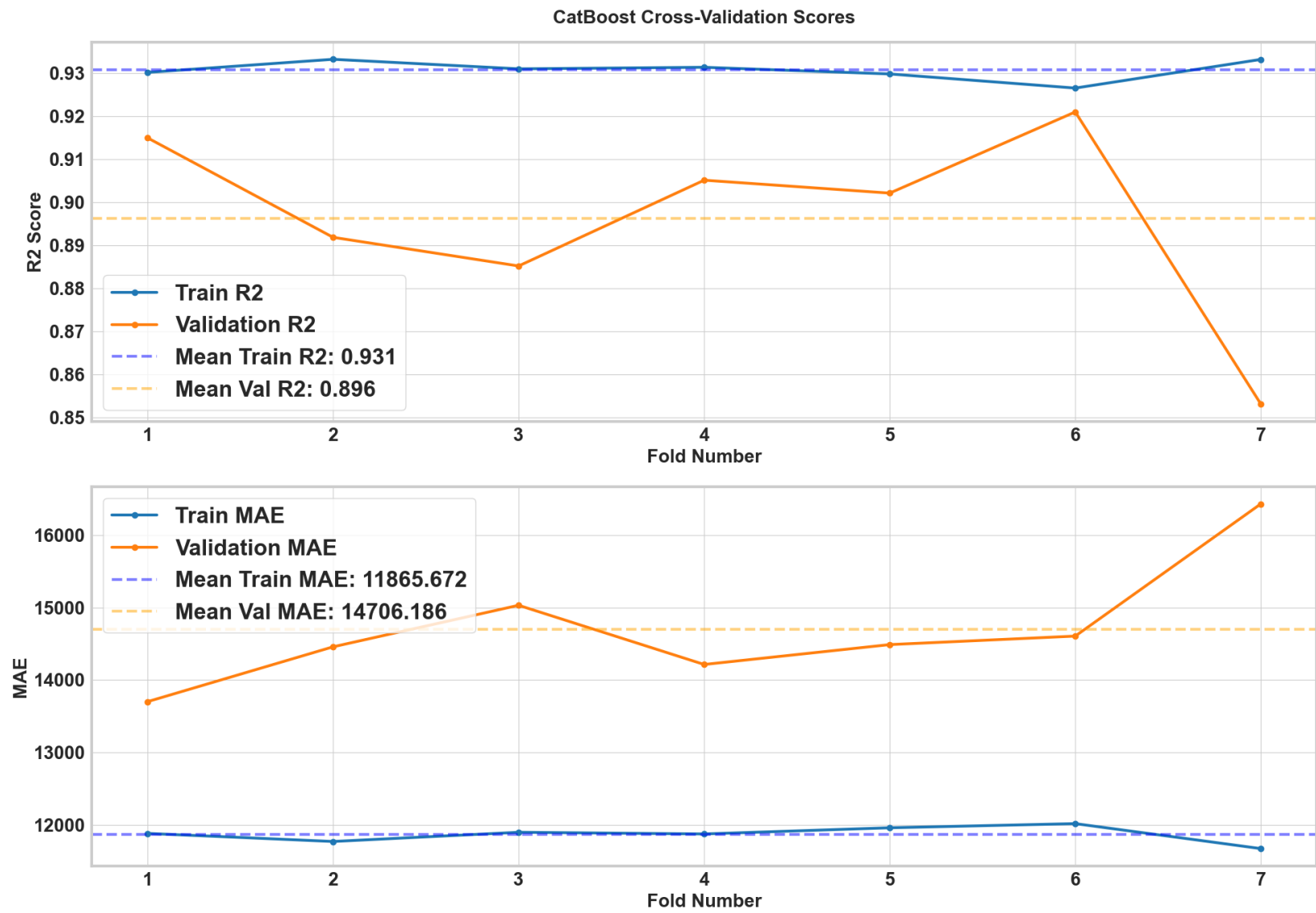
# Cross validation shows slight variation in OOF prediction

LightGBM Cross-Validation Scores

# Cross validation shows slight variation in OOF prediction



CatBoost Cross-Validation Scores

# Cross validation shows slight variation in OOF prediction



AdaBoost Cross-Validation Scores

# All models stabilize for sample sizes above 1500, but they consistently overfit

# LightGBM excels in both predictive power and latency; CatBoost offers slightly enhanced predictive capabilities

## Hyperparameter Tuned Model Performance

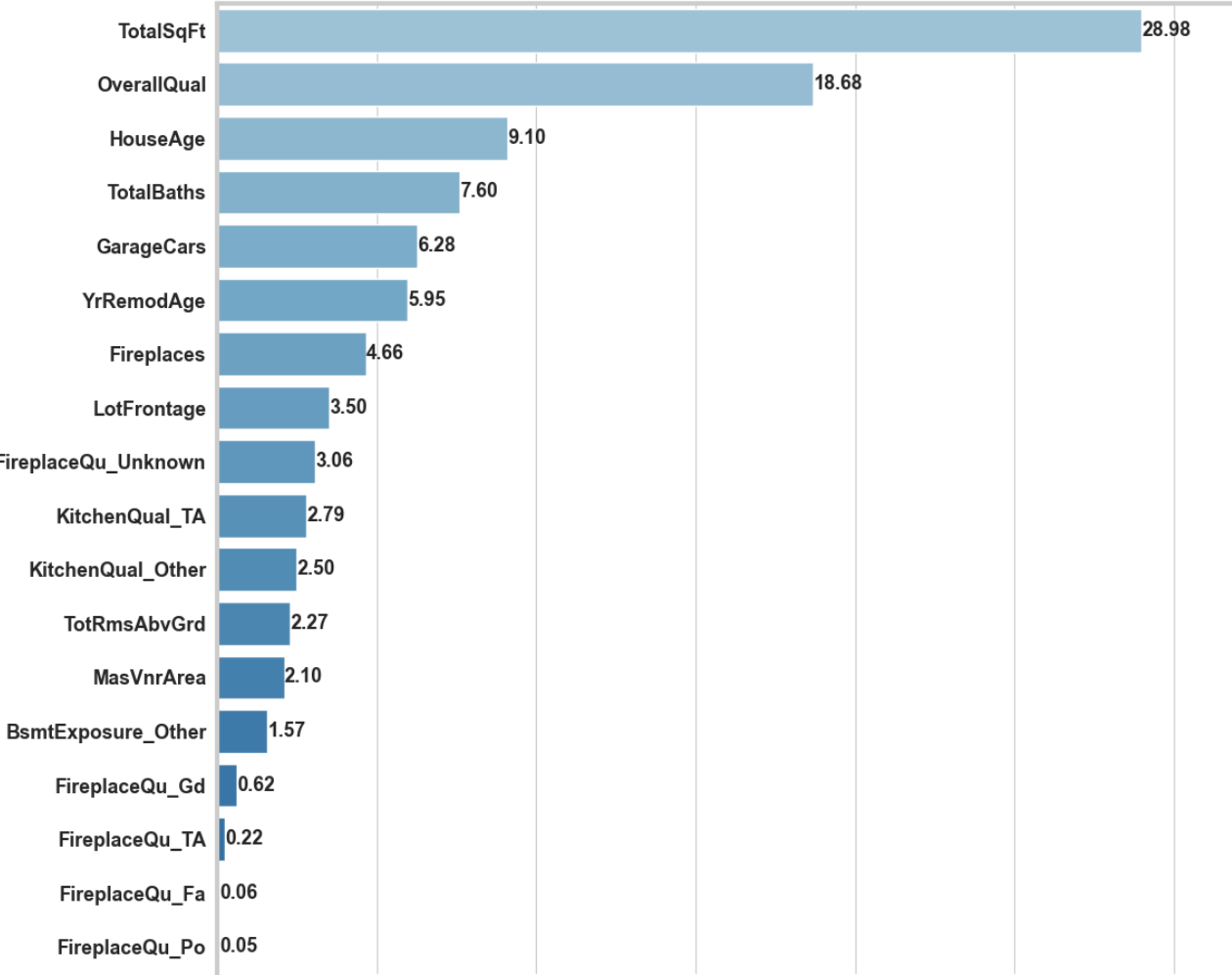| Model | Train R2 | Test R2 | Train MAE | Test MAE | Train MAPE | Test MAPE | Training Time |
|---|---|---|---|---|---|---|---|
| CatBoost | 0.9362 | 0.9386 | 11091.4807 | 12893.8678 | 0.075036 | 0.080996 | 2.41s |
| LightGBM | 0.9371 | 0.9280 | 9783.1770 | 13702.3802 | 0.067831 | 0.087694 | 0.12s |
| RandomForest | 0.9540 | 0.9242 | 10225.5118 | 14146.3613 | 0.067312 | 0.089740 | 1.70s |
| AdaBoost | 0.8378 | 0.8526 | 20584.1489 | 20969.9087 | 0.133344 | 0.125758 | 0.21s |
| Ensemble | 0.9364 | 0.9295 | 12071.4640 | 14215.5670 | NaN | NaN | N/A |

**Latency vs Performance Trade-off**

➢ LightGBM emerges as the most efficient model with just 0.12s training time while maintaining strong performance ($R^2$ 0.928, MAE 13702)

➢ CatBoost achieves the best test accuracy ($R^2$ 0.939) but requires 20x longer training time (2.41s) than LightGBM

➢ RandomForest's higher training $R^2$ (0.954) comes at the cost of longer training time (1.70s), suggesting potential overfitting

➢ The choice between CatBoost and LightGBM would depend on whether the 1% improvement in $R^2$ justifies the 20x increase in training time
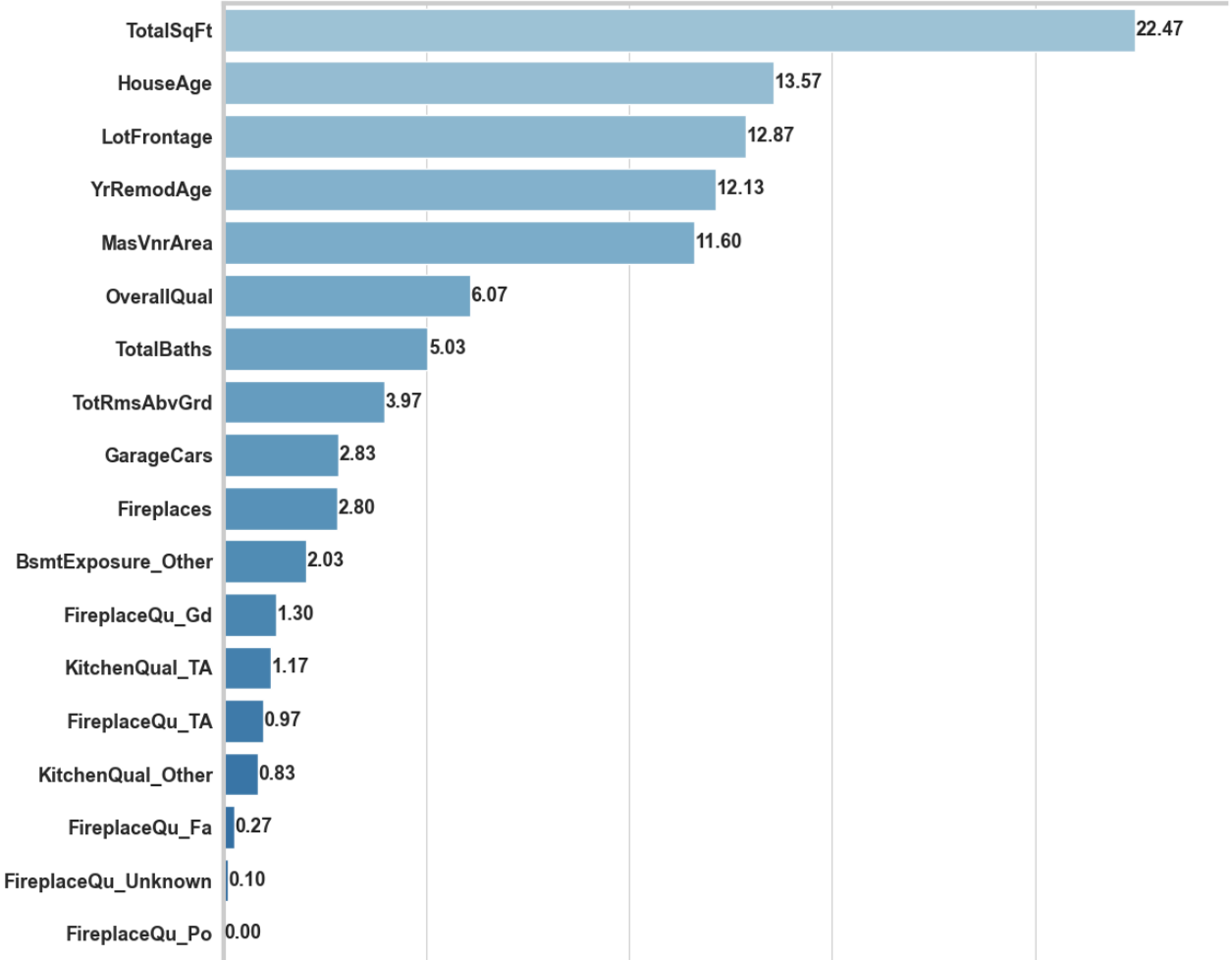
```
best_params = {
    "CatBoost": {
        "colsample_bylevel": 1.0,
        "depth": 5,
        "early_stopping_rounds": 500,
        "eval_metric": "MAE",
        "iterations": 1000,
        "learning_rate": 0.03,
        "min_data_in_leaf": 1,
        "objective": "MAE",
        "subsample": 0.9,
        "verbose": 0
    },
    "LightGBM": {
        "colsample_bytree": 0.9,
        "learning_rate": 0.15,
        "metric": "mae",
        "min_child_samples": 25,
        "min_child_weight": 0.001,
        "n_estimators": 100,
        "num_leaves": 31,
        "objective": "regression_l1",
        "reg_alpha": 0.0,
        "reg_lambda": 0.0,
        "stopping_rounds": 50,
        "subsample": 0.9
    },
    "RandomForest": {
        "max_depth": 9,
        "min_samples_leaf": 2,
        "min_samples_split": 5,
        "n_estimators": 150
    },
    "AdaBoost": {
        "learning_rate": 0.8,
        "loss": "linear",
        "n_estimators": 50
    }
}
```

# Feature importance: How different models tell different stories

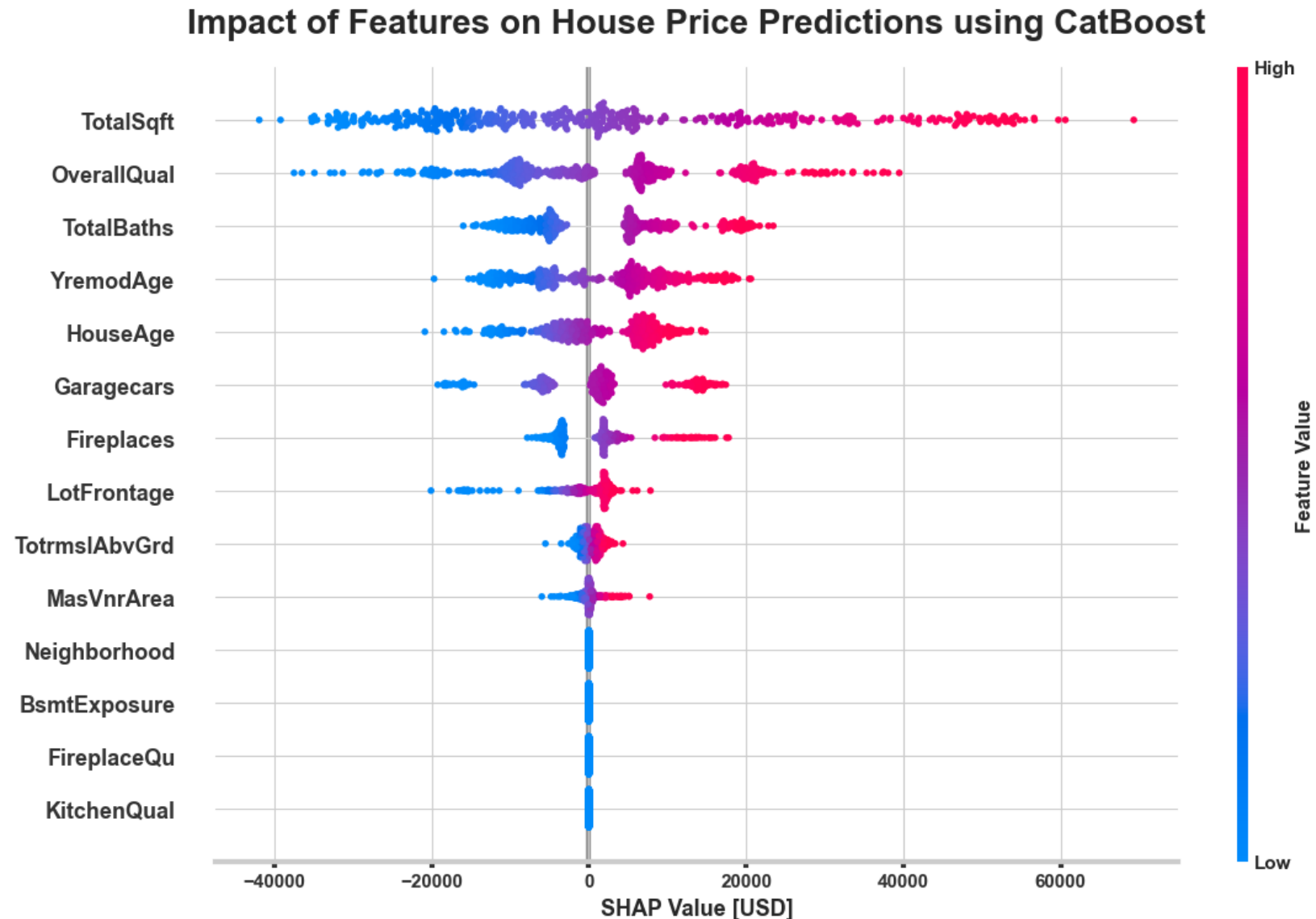### Top 20 Most Important Features for Hypertuned - CatBoost

| Feature | Value |
|---|---|
| TotalSqFt | 28.98 |
| OverallQual | 18.68 |
| HouseAge | 9.10 |
| TotalBaths | 7.60 |
| GarageCars | 6.28 |
| YrRemodAge | 5.95 |
| Fireplaces | 4.66 |
| LotFrontage | 3.50 |
| FireplaceQu_Unknown | 3.06 |
| KitchenQual_TA | 2.79 |
| KitchenQual_Other | 2.50 |
| TotRmsAbvGrd | 2.27 |
| MasVnrArea | 2.10 |
| BsmtExposure_Other | 1.57 |
| FireplaceQu_Gd | 0.62 |
| FireplaceQu_TA | 0.22 |
| FireplaceQu_Fa | 0.06 |
| FireplaceQu_Po | 0.05 |

### Top 20 Most Important Features for Hypertuned - LightGBM

| Feature | Value |
|---|---|
| TotalSqFt | 22.47 |
| HouseAge | 13.57 |
| LotFrontage | 12.87 |
| YrRemodAge | 12.13 |
| MasVnrArea | 11.60 |
| OverallQual | 6.07 |
| TotalBaths | 5.03 |
| TotRmsAbvGrd | 3.97 |
| GarageCars | 2.83 |
| Fireplaces | 2.80 |
| BsmtExposure_Other | 2.03 |
| FireplaceQu_Gd | 1.30 |
| KitchenQual_TA | 1.17 |
| FireplaceQu_TA | 0.97 |
| KitchenQual_Other | 0.83 |
| FireplaceQu_Fa | 0.27 |
| FireplaceQu_Unknown | 0.10 |
| FireplaceQu_Po | 0.00 |

**CatBoost** emphasizes **quality metrics**: OverallQual, TotalBaths, and KitchenQual rank higher. **LightGBM** gives more weight to **physical attributes**: LotFrontage and MasVnrArea show much higher importance

# $$ value impact: What each home feature adds

- ➢ Total Square Footage (TotalSqft) has the largest impact range (~60k USD) and shows consistently positive influence for larger values

- ➢ Overall Quality (OverallQual) is the second most influential feature, with higher quality scores strongly driving up prices

- ➢ Neighborhood and quality-related categorical features (KitchenQual, FireplaceQu, BsmtExposure) show clustered impacts



Impact of Features on House Price Predictions using CatBoost

# Potential Directions

➢ **LLMs in Proptech**: Utilizing LLMs to scan and extract valuable information from property and legal documents, enhancing property transaction efficiency.

➢ **Foreclosure Predictions**: Implementing AI models to forecast foreclosure risks and assess buyer preparedness after listing a property for sale.

➢ **Image Analysis for Property Assessment**: Using AI for image analysis to detect property damages, aiding in accurate property valuation.

➢ **Proactive Real Estate Services**: Collaborating with banks and mortgage lenders to create databases tracking mortgage defaults, offering proactive property management.

➢ **AI/ML as SaaS**: Providing AI and ML technologies as a service to the real estate sector, enabling unprecedented insights and operational efficiency.

# Acknowledgements

I would like to extend my heartfelt gratitude to everyone
who provided invaluable assistance during this project.
- ❖ Vivian S. Zhang
- ❖ Cole Ingraham
- ❖ My cohorts: James Seykot, Margaret Bowers, Oreste
Rukundo, Amiyo Chattarjee