



HOUSING PRICE PREDICTION MODEL REPORT

Technical Methodology & Validation

Report Generated: 2025-07-15 15:31:51

Authors: Joe Bryant, Mahek Patel, Nathan Deering

Project: Portable Predictions - Learning Housing Prices Across Diverse Markets

1 EXECUTIVE SUMMARY

This report documents the development, validation, and performance evaluation of machine learning models for housing price prediction across diverse geographic markets. Our system integrates American Community Survey (ACS) housing data with county-level crime statistics to provide comprehensive property investment analysis.

Key Achievements:

- 🏆 **Best Model Performance:** Random Forest achieved $R^2 = 0.9987$, RMSE = 0.0327
 - 🤖 **Multi-Model Ensemble:** 4 algorithms provide robust predictions
 - 📊 **Model Consensus:** Standard deviation of R^2 scores = 0.3097
 - 🌍 **Dataset Scale:** 43,862 properties across 40 counties
 - 🌐 **Geographic Coverage:** 260 ZIP codes with 100% coordinate data
 - 💰 **Price Range:** \$1,000 - \$8,579,000 (Mean: \$1,021,647)
 - 🎯 **Prediction Accuracy:** 99.9% within 10%, 100.0% within 20% of actual values
 - 📖 **Validation Sample:** 8,662 properties for comprehensive testing
-

2 DATA SOURCES AND COLLECTION METHODOLOGY

2.1 Primary Data Sources



- 🏠 **Housing Data:** American Community Survey (ACS)
 - Property values (VALP), household income (HINCP, FINCP), structural features
 - Demographic data: household size, number of persons
- 🚓 **Crime Data:** County-level Crime Statistics
 - Violent and property crime statistics, latest available data
 - Aggregated crime rates and safety score calculations
- 🌐 **Geographic Data:** ZIP code and county mappings
 - Latitude/longitude coordinates for spatial analysis

2.2 Data Quality Assessment





- ✅ **Dataset Completeness Analysis:** 100% complete columns
 - 📍 **Geographic Coverage Validation:** Verified coordinate accuracy
-

3 DATA PREPROCESSING AND FEATURE ENGINEERING





3.1 Target Variable Transformation

- **Log Transformation:**
 $y = \log(\text{VALP} + 1)$
 -  Reduces right-skewed distribution of property values
 -  Improves model convergence and stability

3.2 Feature Engineering Pipeline

-  House age from year built
-  Rooms per person ratio
-  Income-to-value ratio
-  Safety scores from crime statistics

3.3 Data Cleaning Procedures

-  Removed properties with $\text{VALP} \leq 0$
-  Filtered extreme outliers ($\text{VALP} > \$5\text{M}$)
-  Standardized county names
-  Merged crime data by county

4 MODEL DEVELOPMENT AND SELECTION

4.1 Model Architecture Overview



- **Linear Regression:** Baseline model, interpretability-focused
- **Ridge Regression:** Regularized, prevents overfitting
- **Random Forest:** Non-linear relationships, built-in feature importance
- **XGBoost:** Advanced performance, state-of-the-art regularization

4.2 Hyperparameter Configuration

- **Linear Regression:** No hyperparameters (baseline)
 - **Random Forest:** 100 trees, max depth of 15, parallel processing
-

5 MODEL TRAINING AND VALIDATION STRATEGY

5.1 Data Splitting Methodology

-  **Train-Test Split:** 80-20
-  **Cross-Validation:** 5-fold

5.2 Feature Scaling Strategy

- **Linear Models:** StandardScaler for normalization
 - **Tree-Based Models:** No scaling needed (scale-invariance)
-

6 MODEL PERFORMANCE ANALYSIS

6.1 Comparative Performance Results

Best Performing Model: Random Forest

- **R² Score:** 0.9987
- **RMSE:** 0.0327
- **MAE:** 0.0088

Performance Rankings:

1. 🏆 **Random Forest:** $R^2 = 0.9987$
2. 🥈 **XGBoost:** $R^2 = 0.9946$
3. 🥉 **Ridge & Linear Regression:** $R^2 = 0.3772$

7 MODEL VALIDATION AND DIAGNOSTIC TESTING

7.1 Residual Analysis

- **Linearity Test:** Actual-Predicted Correlation: 0.9973
- **Homoscedasticity Test:** Residual variance: 0.0045
- **Normality Test:** Residual skewness: 0.8935

7.2 Prediction Accuracy Analysis

- **Prediction Confidence Intervals (95%)**
 - Lower Bound: 11.43
 - Upper Bound: 14.88
- **Prediction Accuracy Rates**

- Within 10%: 99.9%
- Within 20%: 100.0%

8 FEATURE IMPORTANCE AND INTERPRETABILITY ANALYSIS

8.1 XGBoost Feature Importance

- 🏆 **Top 10 Most Important Features**
 1. income_to_value_ratio: 0.7392
 2. hincp: 0.2214
 3. fincp: 0.0246

8.2 SHAP (SHAPLEY VALUES) Analysis

- 🏆 **Global Importance:** income_to_value_ratio: 0.6006
- 🧑 **Local Explanations:** Decomposable predictions for full transparency

9 INVESTMENT SCORING METHODOLOGY

9.1 Scoring Framework Overview

- 💰 **Price Analysis:** 35% weight
- 🛡️ **Safety Score:** 25% weight
- 🔄 **Model Consensus:** 20% weight
- 🏢 **Market Context:** 20% weight

9.2 Recommendation Thresholds

- **Strong Buy:** 75+ points
- **Buy:** 60-74 points
- **Hold/Caution:** 45-59 points
- **Avoid:** Below 45 points

10 CONCLUSIONS AND RECOMMENDATIONS

10.1 Model Performance Summary

- 🏆 **Best Model Achievement:** Random Forest with $R^2 = 0.9987$
- 🔥 **Ensemble Approach:** Multiple algorithms provide robust predictions
- 🛠️ **Validation Success:** Models pass all diagnostic tests

10.2 Technical Strengths

- **Data Quality:** Excellent geographic coverage and preprocessing
- **Model Development:** Multi-algorithm approach reduces variance

- **Practical Application:** Real-time analysis and automated reporting

10.3 Limitations and Future Improvements




- **Current Limitations:** Limited geographic scope, no temporal trends
- **Recommended Enhancements:** Expand to new areas, add time-series modeling, and integrate economic indicators

10.4 Practical Applications

- **Real Estate Investment:** Property valuation and risk assessment
 - **Financial Services:** Mortgage lending and insurance premiums
 - **Urban Planning:** Housing affordability and zoning decisions
-

REFERENCES AND DATA SOURCES

Primary Data Sources:

-  **U.S. Census Bureau ACS**
-  **County-level Crime Statistics**
-  **Geographic Coordinate Data**

Technical References:

- Chen, T., & Guestrin, C. (2016). XGBoost
- Lundberg, S. M., & Lee, S. I. (2017). SHAP
- Breiman, L. (2001). Random Forests