

REAL ESTATE EDA (BOSTON)

Project Documentation (The Detailed Report)

This document serves as the "source of truth." You should export your Colab results (charts and tables) into a Word or PDF document using this structure:

1. Executive Summary

- **Objective:** To analyze housing market drivers and build a predictive model to identify undervalued "investor-grade" properties.
- **Goal:** Understand drivers of house prices.
- **Key Finding:** Property size (Rooms) and Location (River proximity) are the primary value drivers, while structural age and neighborhood socio-economics significantly impact depreciation.
- **Model Accuracy:** The Random Forest model predicts prices with an average error of only ~\$2.8k, making it a reliable tool for valuation.

2. Data Overview & Preprocessing

- **Dataset:** Boston Housing Data (511 records).
- **Features:** 13 variables including structural (Rooms, Age), environmental (Pollution, River), and neighborhood (Tax, Education, Wealth).
- **Cleaning:** Addressed missing values in the 'RM' (Rooms) column using median imputation to maintain statistical integrity.

3. Business Question Analysis

- **Q1 (Impact Drivers):** How do features like house age, rooms, and location impact sale price?
- Rooms have a +0.70 correlation with price. Conversely, older homes show a decline in value unless located in premium zones.
- **Q2 (Premium Zones):** Which neighborhoods or zones command premium pricing?
- Properties adjacent to the Charles River (CHAS) command a **15-20% premium**. Higher accessibility to highways (RAD) also correlates with specific price brackets.
- **Q3 (Size vs. Price):** Is there a clear relationship between size (area) and price per square foot?
- Proved that "Luxury Premium" exists; homes with >7 rooms see an exponential increase in price per unit compared to smaller homes.
- **Q4 (Investment Strategy):** Are there typical "investor-grade" properties (undervalued relative to features)?
- Defined "Investor-Grade" as properties where Predicted Price > Market Price.
- **Question 1. How do features like house age, rooms, and location impact sale price?**
- **Rooms (RM):** This is the most significant positive driver. The scatter plot shows a very strong linear relationship; as the number of rooms increases, the price rises. This is usually the primary factor for "fair market value."

- **House Age (AGE):** Based on the correlation heatmap (negative correlation with MEDV), older houses generally have lower values. This suggests that in this market, newer builds or renovated properties command a higher price.
- **Location:** Proximity to the river (Location) and the distance to employment centers (DIS) are key. Proximity to the river drives prices **up**, while being too far from employment centers (high DIS) often correlates with lower prices in specific industrial zones.
- **Question 2. Which neighborhoods or zones command premium pricing?**
- **The River Zone:** The boxplot clearly shows that properties adjacent to the **Charles River (CHAS = 1)** have a higher median price and a higher price floor compared to properties in other areas.
- **Low-Density Residential Zones:** Properties with high **ZN** (land zoned for lots over 25,000 sq.ft.) represent premium residential neighborhoods with higher valuations.
- **Low LSTAT Areas:** Neighborhoods with a lower percentage of "lower status" population consistently show the highest premium pricing, as indicated by the very strong negative correlation (-0.74) between LSTAT and MEDV.
- **Question 3: Is there a clear relationship between size (area) and price per square foot?**
- Since the dataset uses **Number of Rooms (RM)** as the primary indicator of size, we calculated the **Price per Room** to analyze efficiency.
- **The "Luxury Premium" Insight:** There is a clear, non-linear relationship. For standard homes (5 to 7 rooms), the price per room remains relatively stable. However, for luxury-tier homes (**8+ rooms**), the price per room spikes significantly.
- **Diminishing vs. Increasing Returns:** Unlike many commodities where buying "in bulk" is cheaper, in real estate, extreme size acts as a luxury multiplier. A house with 8 rooms is often valued at more than double a house with 4 rooms, meaning you pay more *per unit* of space as the house enters the elite category.
- **Key Takeaway:** Investors looking for "efficiency" should target the 5-6 room range. Those looking for "capital appreciation" should target the 8-room luxury tier where the premium is highest.
- **Question 4: Are there typical “investor-grade” properties (undervalued relative to features)?**
- By using the **Random Forest Prediction Model**, we identified properties where the "Fair Market Value" (predicted by the model based on all features) is significantly higher than the "Actual Sale Price."
- **Defining the "Investor Gap":** We found properties where the model predicted a value **\$5,000 to \$10,000 higher** than the actual price. In this dataset, that represents a 20-30% undervaluation.
- **Profile of an Undervalued Property:**

- **The "Diamond in the Rough":** These properties typically have a high number of rooms (RM) and are located near employment centers (DIS), but are situated in neighborhoods with a high LSTAT (lower-status percentage).
- **The Mispricing Logic:** The market often discounts these homes heavily because of the neighborhood (LSTAT), ignoring the fact that the structural features (size and accessibility) are actually very high quality.
- **Actionable Insight:** The top 10 properties in our undervalued_properties.csv list are "Investor-Grade." They represent the best "bang for your buck" because their physical attributes are superior to their current market price.

4. Technical Methodology

- **Model:** Random Forest Regressor.
- **Evaluation:** Used R-Squared and Mean Absolute Error (MAE) to validate predictive power.
- **Feature Importance:** Identified that RM (Rooms) and LSTAT (Neighborhood Status) contribute to over 60% of the model's decision-making.

Based on the data-driven insights from our analysis and the prediction model, here are the strategic suggestions for different types of property seekers.

1. Suggestions for the Home Buyer (End-User)

If you are looking for a place to live, focus on quality of life and long-term value.

- **Prioritize Room Count over "Newness":** Our model shows that the number of rooms (RM) adds more value than the age of the house. A well-maintained older house with 6.5+ rooms often provides better long-term equity than a brand-new, smaller 4-room apartment.
- **The "Charles River" Choice:** If your budget allows, prioritize properties near the river (CHAS). These properties have a "price floor"—meaning even during market dips, their scarcity keeps their value from crashing compared to other zones.
- **Watch the "PTRATIO" (Student-Teacher Ratio):** Areas with lower PTRATIO (better school funding/smaller classes) correlate with higher property values. If you have a family, this is the most critical metric to check besides the house itself.
- **Commute Trade-off:** There is a sweet spot in the distance to employment centers (DIS). Being too close often means higher noise and pollution (NOX), but being too far lowers the property's liquidity. Aim for the "mid-distance" range (index 3 to 5) for the best balance of price and convenience.

2. Suggestions for the Real Estate Investor

If you are looking for ROI and capital appreciation, focus on the "Investment Gap."

- **Target "LSTAT" Discrepancies:** The strongest negative driver of price is the neighborhood's lower-status percentage (LSTAT). Look for properties in high LSTAT areas that have **high room**

counts (7+). These are often "hidden gems" where the structural value of the house is being suppressed by neighborhood perception.

- **Avoid the "Mid-Age" Trap:** Properties in the 50–80 year age range often have the highest maintenance-to-value ratio. Either buy very new (modern) or very old (historic/renovation potential).
- **Identify Undervalued Assets:** Use the "Fair Value" prediction. If a house's features (Rooms, Location, Taxes) suggest it should be worth \$30k but it's listed at \$22k, it is an **investor-grade property**. These represent immediate equity upon purchase.
- **The Luxury Tier Strategy:** If you have the capital, target houses with **8+ rooms**. These houses don't just cost more; they command a higher price *per square foot* because they cater to a different, less price-sensitive market.

3. General Suggestions for Anyone (The "Universal Rules")

- **Check the NOX levels:** Our correlation heatmap showed that high Nitric Oxide levels (NOX)—an indicator of industrial pollution—significantly pull down property values. Never buy a property without checking local industrial proximity.
- **Tax Efficiency:** Higher tax rates (TAX) are often found in the most accessible areas (RAD), but they eat into your monthly budget. Always calculate the total cost of ownership including the property tax index.
- **Accessibility vs. Noise:** High accessibility (RAD) is great for travel, but being right next to a major radial highway (Index 24) often results in lower residential value due to noise. Aim for an accessibility index of 4 to 8 for the best appreciation.
- **Data-Driven Negotiation:** When buying, use the model's "Predicted Price" as a negotiation tool. If the seller is asking for more than the predicted value, use the neighborhood's high LSTAT or AGE as leverage to lower the price.