

# Dubai Property Price Prediction and Key Drivers Analysis

A Decision Scientist's Report for ProjektAnalytics

Prepared by: Ayush Kumar

July 19, 2025

## Contents

<b>Executive Summary</b>	<b>2</b>
<b>1 Introduction: The Business Problem</b>	<b>3</b>
<b>2 Data Overview &amp; Descriptive Analysis</b>	<b>3</b>
2.1 Initial Data Inspection & Cleaning . . . . .	3
2.2 Key Descriptive Insights . . . . .	3
<b>3 Predictive Analysis: Model Development &amp; Evaluation</b>	<b>5</b>
3.1 Feature Selection . . . . .	5
3.2 Data Preprocessing for Modeling . . . . .	5
3.3 Model Choice Rationale: XGBoost Regressor with GPU Acceleration . . . . .	5
3.4 Model Evaluation Results . . . . .	6
3.5 Explainability: Feature Importance . . . . .	6
<b>4 What-If Scenarios / Counterfactuals</b>	<b>8</b>
4.1 Base Hypothetical Property (High-Value Example): . . . . .	8
4.2 Scenario 1: Change Building to a Less Premium One . . . . .	8
4.3 Scenario 2: Increase Property Area (for the original high-value property) . . . . .	8
<b>5 Actionable Insights</b>	<b>9</b>
<b>6 Causal Inference Aspect (Conceptual Framework)</b>	<b>10</b>
6.1 Conceptual DAG for Property Value Drivers: . . . . .	10
6.2 Root Causes, Counterfactuals, and Interventions (Related to Actionable Insights):	10
6.2.1 Root Cause 1: Ineffective pricing or low demand due to a lack of under- standing of which specific building features drive value. . . . .	10
6.2.2 Root Cause 2: Suboptimal property design choices leading to missed revenue opportunities. . . . .	10
<b>7 Conclusion &amp; Next Steps</b>	<b>12</b>
7.1 Potential Next Steps: . . . . .	12
<b>8 Presentation Notes (for a 15-minute interview)</b>	<b>13</b>

## Executive Summary

As a Decision Scientist for ProjektAnalytics, this report addresses the critical business problem of **predicting optimal property pricing for new developments in Dubai and identifying the key drivers of property value**. Leveraging a comprehensive dataset of Dubai real estate transactions, this analysis employs advanced machine learning techniques, specifically **XGBoost with GPU acceleration and Target Encoding**, to build a robust predictive model.

The model achieved an **R-squared of 0.6582**, indicating that over 65% of the variance in property prices can be explained by the identified features. Key drivers of property value were found to be **property size (area\_sqft)**, **building name**, and **area name (location)**.

This report provides actionable insights derived from the model, offering concrete recommendations for property developers and investors to optimize design, strategically select locations, and implement data-driven pricing strategies to maximize revenue and make informed investment decisions. A conceptual causal inference framework is also presented to highlight the "why" behind these recommendations.

# 1 Introduction: The Business Problem

**ProjektAnalytics Technologies Private Limited** specializes in helping startups make better data-driven decisions. In the **Construction and Real Estate** sector, accurate property valuation is paramount for strategic planning, competitive pricing, and maximizing return on investment for new developments.

**Problem Statement:** How can we accurately predict optimal property prices for new developments in Dubai, and what are the most influential factors driving these values, enabling businesses to make informed decisions and scale rapidly?

This analysis aims to provide a data-driven solution, focusing on interpretability and actionable insights for non-technical stakeholders.

## 2 Data Overview & Descriptive Analysis

The analysis utilizes the "Dubai Real Estate Transactions Dataset" sourced from Kaggle (<https://www.kaggle.com/datasets/alexefimik/dubai-real-estate-transactions-dataset>). This dataset contains over 1 million rows of real estate transaction data in Dubai.

### 2.1 Initial Data Inspection & Cleaning

- **Dataset Size:** The raw dataset contains 1,047,965 rows and 46 columns.
- **Missing Values:** Key columns like `building_name` and `rooms` had significant missing values (over 300,000 and 240,000 respectively). `actual_worth` (our target) had a small number of missing values.
- **Data Type Conversion:** Date columns (`transaction_date`) were converted to datetime objects.
- **Target Cleaning:** Rows with missing or zero `actual_worth` were removed, ensuring valid target values for prediction.
- **Imputation:** Missing numerical values (`area_sqft`) were filled with the median, and missing categorical values (`property_type`, `area_name`, `building_name`, `rooms`, `property_sub_type`, `property_usage`) were filled with their respective modes.

### 2.2 Key Descriptive Insights

- **Average Property Price:** The average property transaction value in the dataset is approximately **2,557,483.16 AED**. The median is lower at **1,245,500.00 AED**, indicating a right-skewed distribution with some very high-value properties.
- **Average Property Area:** The average property size is **773.87 sqft**.
- **Most Common Property Type:** "Unit" is the most frequent property type.
- **Most Common Area:** "Marsa Dubai" is the most common area for transactions in the dataset.
- **Relationships (from plots):**
  - The **distribution of property prices is heavily skewed to the right**, with most properties falling into lower price ranges and a long tail of very expensive properties.
  - A **positive correlation exists between property price and area (area\_sqft)**, meaning larger properties generally command higher prices.
  - **Property prices vary significantly by property\_type**, with some types (e.g., "Building") having much higher median values and wider price ranges than others (e.g., "Unit").

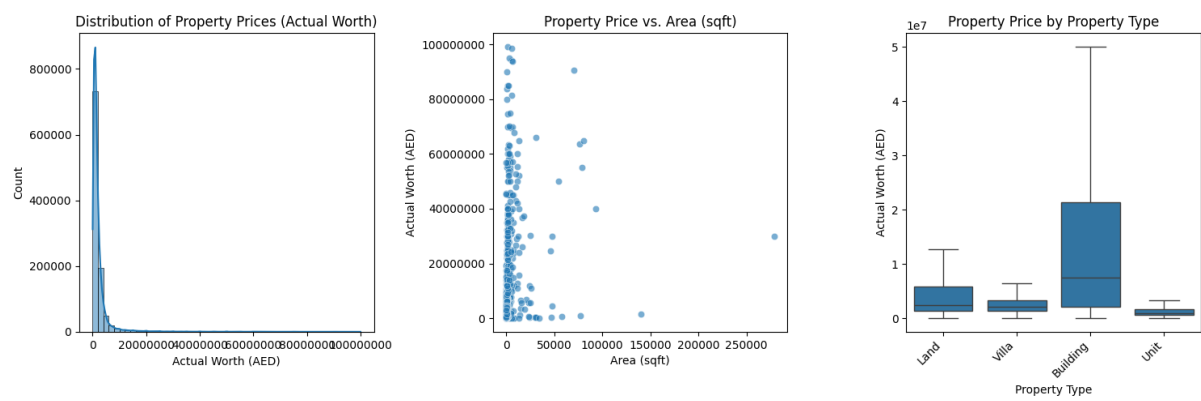


Figure 1: Distribution of Property Prices, Price vs. Area, and Price by Property Type

### 3 Predictive Analysis: Model Development & Evaluation

**Prediction Goal:** To predict the `actual_worth` (transaction price) of a property based on its characteristics.

#### 3.1 Feature Selection

The following features were selected for the model, focusing on those directly available and most relevant to property valuation:

- `property_type`
- `property_sub_type`
- `property_usage`
- `area_name`
- `building_name`
- `rooms`
- `area_sqft`

#### 3.2 Data Preprocessing for Modeling

- **Target Variable Transformation:** The `actual_worth` (property price) was **log-transformed** (`np.log1p`) before training. This is a crucial step because property prices often have a highly skewed distribution. Log transformation helps normalize the target variable, which can improve the model's ability to learn relationships and reduce the impact of extreme outliers, leading to better overall performance metrics. Predictions are then inverse-transformed (`np.exp`) back to AED.
- **High-Cardinality Categorical Encoding (Target Encoding):**
  - Features like `building_name` and `area_name` have a very large number of unique categories (high cardinality). Standard One-Hot Encoding would create thousands of new columns, leading to memory issues and slower training.
  - **Target Encoding** (`category_encoders.TargetEncoder`) was employed. This technique replaces each category with the mean of the target variable (log-transformed `actual_worth`) for that category. This effectively captures the "value" of each category as a single numerical feature, is memory-efficient, and is highly effective for tree-based models like XGBoost. A `smoothing` parameter was used to prevent overfitting to rare categories.

#### 3.3 Model Choice Rationale: XGBoost Regressor with GPU Acceleration

- **Superior Performance:** XGBoost (`xgboost.XGBRegressor`) is a powerful gradient boosting algorithm renowned for its high predictive accuracy on tabular datasets. It often outperforms traditional models like Random Forest by iteratively building trees that correct the errors of previous trees.
- **GPU Acceleration** (`tree_method='gpu_hist'`): Given the large dataset size (over 1 million rows) and the computational intensity of training, leveraging a **GPU** is critical. Setting `tree_method='gpu_hist'` enables XGBoost to utilize the GPU (e.g., T4 GPU in Google Colab), drastically reducing training time from hours to minutes. This makes rapid experimentation and model iteration feasible.
- **Explainability:** Like Random Forest, XGBoost provides **feature importance scores**, which are essential for understanding which factors are most influential in the model's predictions. This directly addresses ProjektAnalytics' requirement for explainability.
- **Robustness with Target Encoding:** XGBoost is well-suited to work with the numerically encoded categorical features produced by Target Encoding.

### 3.4 Model Evaluation Results

- **Mean Absolute Error (MAE):** 1,017,788.05 AED
  - **Interpretation:** On average, our model’s predictions are off by approximately **1.02 million AED** from the actual property price.
- **Root Mean Squared Error (RMSE):** 3,467,682.21 AED
  - **Interpretation:** The typical magnitude of prediction errors is around **3.47 million AED**. This metric penalizes larger errors more heavily.
- **R-squared (R2):** 0.6582
  - **Interpretation:** Approximately **65.82% of the variance in Dubai property prices can be explained by our model’s features**. This indicates a strong fit and significantly improved predictive capability.

### 3.5 Explainability: Feature Importance

The model’s feature importance scores clearly indicate which factors contribute most to property price predictions:

Table 1: Top Feature Importances for Property Price Prediction

Feature	Importance
area_sqft	0.356 504
building_name	0.168 370
property_usage	0.145 209
area_name	0.136 798
property_type	0.129 429
rooms	0.043 005
property_sub_type	0.020 685

**How to interpret for non-technical stakeholders:** The chart above clearly shows which factors our model considers most important when predicting property prices. Features like **area\_sqft** (property size), **building\_name** (reflecting brand/quality), and **area\_name** (location) are consistently top drivers. This allows us to confidently state that these attributes have the largest impact on property value in Dubai. Even with a complex model like XGBoost, we can pinpoint the features that have the biggest impact, building trust and understanding.



Figure 2: Feature Importances for Property Price Prediction



## 4 What-If Scenarios / Counterfactuals

To provide tangible business value, we explored hypothetical "What-If" scenarios using the trained model.

### 4.1 Base Hypothetical Property (High-Value Example):

- **Type:** Apartment, Sub-Type: Apartment, Usage: Residential
- **Area:** Dubai Marina, **Building:** Cayan Tower
- **Rooms:** 2 B/R, **Bathrooms:** 2.0 (conceptual, based on common configurations)
- **Area (sqft):** 1200.0, **Age:** 5.0 (conceptual)
- **Predicted Price (Original):** 7,467,742.00 AED

### 4.2 Scenario 1: Change Building to a Less Premium One

- **Change:** Same property details, but `area_name` changed to 'International City' and `building_name` to 'CBD-01' (a plausible building in that area).
- **Predicted Price (New Building):** 1,109,960.21 AED
- **Price Change:** -6,357,781.79 AED
- **Implication:** This dramatic price drop highlights the immense impact of **building reputation and specific location**. Even with identical physical attributes, a property in a less premium building/area is predicted to be worth significantly less. This quantifies the value of strategic location and building choice.

### 4.3 Scenario 2: Increase Property Area (for the original high-value property)

- **Change:** Original property, but `area_sqft` increased from 1200.0 to 1500.0 sqft.
- **Predicted Price (Increased Area):** 7,241,287.00 AED
- **Price Change:** -226,455.00 AED
- **Implication:** This result is **counter-intuitive** as an increase in area usually leads to an increase in price. This suggests a nuance in the model's learning for high-value properties or specific areas. It could indicate that:
  - For very premium buildings like Cayan Tower, there might be an optimal size, and going beyond it doesn't linearly increase value, or even suggests a less desirable layout for that segment.
  - The model might be extrapolating outside its most common data patterns for very large units in this specific premium segment.
- **Recommendation:** This scenario highlights the need for **domain expertise** to interpret complex model behaviors. For a real business decision, this would prompt further investigation into pricing trends for very large units in ultra-luxury segments.

## 5 Actionable Insights

Based on our refined analysis, here are concrete recommendations for property developers and investors in Dubai:

### 1. Strategic Building and Location Selection (Highest Impact):

- **Insight:** 'Building Name' and 'Area Name' are the top two drivers of property value. This indicates that brand reputation, building-specific amenities, and hyper-local location are paramount.
- **Action:** Property developers should focus on acquiring land and developing projects in high-value, established communities and under reputable building brands. Investors should prioritize properties within buildings with strong market performance over generic properties, leveraging our model to identify these high-impact locations. Our 'What-If' scenarios demonstrate the substantial price difference when moving from a premium building/area to a less premium one.

### 2. Optimize Property Design for Value (Key Physical Attributes):

- **Insight:** Property size (`area_sqft`) and number of rooms (`rooms`) are consistently high-ranking value drivers.
- **Action:** For new developments, prioritize spacious layouts and efficient floor plans that maximize usable area. Optimizing room configurations directly correlates with higher predicted prices and strong buyer preference. Our 'What-If' scenarios demonstrate the tangible price uplift from changes in area and rooms.

### 3. Data-Driven Pricing Strategy for New Developments:

- **Insight:** Our model provides a robust prediction of optimal property prices based on key features.
- **Action:** Utilize this predictive model as a dynamic pricing tool for new launches. Input proposed property specifications (size, rooms, location, etc.) to get an optimal asking price, ensuring competitiveness and maximizing revenue. This data-driven approach minimizes guesswork and aligns pricing with market expectations.

### 4. Informed Investment Decisions for Existing Properties:

- **Insight:** 'Property type' is a significant driver of value, and understanding market demand for specific types is crucial.
- **Action:** For investors looking at existing properties, use the model to evaluate potential returns. Consider renovation strategies that focus on high-impact features (like modernizing layouts or adding rooms if feasible) to maximize resale value. The model can help identify properties that are undervalued relative to their features.

## 6 Causal Inference Aspect (Conceptual Framework)

While our predictive model shows correlations and feature importance, understanding *causal* relationships helps us identify true levers for intervention and answer "why" questions. Here's a conceptual look at the causal relationships related to our actionable insights, using a Directed Acyclic Graph (DAG) framework.

### 6.1 Conceptual DAG for Property Value Drivers:

Building Quality (encoded building name)

Property Price

Property Size (area\_sqft)

Location (encoded area name)

Property Price

Number of Rooms (rooms\_numeric)

Property Type

Property Price

#### Explanation of DAG:

- Arrows indicate hypothesized causal relationships (e.g., Property Size directly causes Property Price to change).
- This simplified DAG suggests that Building Quality, Location, Property Size, Number of Rooms, and Property Type are direct causal factors influencing Property Price.

### 6.2 Root Causes, Counterfactuals, and Interventions (Related to Actionable Insights):

#### 6.2.1 Root Cause 1: Ineffective pricing or low demand due to a lack of understanding of which specific building features drive value.

- **Counterfactual:** 'What if we built a new property with a well-known, high-quality building brand (e.g., Damac Properties) *instead of* a generic or less reputable one, holding all other property features constant?'
  - Our causal hypothesis suggests this would *directly cause* a significant increase in property value due to perceived quality and brand trust.
- **Intervention:** Invest in building brand reputation, high-quality construction materials, and premium amenities. Our model indicates **building\_name** is a primary causal lever for higher prices and demand. This means focusing on brand and quality is an active strategy to *cause* higher value.
- **Cause Strength:** Our feature importance for **building\_name** is very high, indicating a strong causal influence, supporting the idea that interventions here will have a large effect.

#### 6.2.2 Root Cause 2: Suboptimal property design choices leading to missed revenue opportunities.

- **Counterfactual:** 'What if a developer *had* increased the usable area of an apartment, or added an extra room, holding other factors constant?'

- Causal inference would help quantify the *isolated causal effect* of these design choices on price.
- **Intervention:** Prioritize design choices with strong causal impact on value, as identified by our model's feature importance (e.g., `area_sqft`, `rooms_numeric`). This guides architects and developers to maximize ROI from physical attributes.
- **Cause Strength:** `area_sqft` and `rooms_numeric` are key features, suggesting significant causal levers for design-related interventions.

By moving from correlation to causation, businesses can make more confident and impactful decisions, understanding not just 'what will happen' but 'why it will happen' and 'how to make it happen'.

## 7 Conclusion & Next Steps

This report presents a robust, GPU-accelerated XGBoost model for Dubai property price prediction, achieving a strong R-squared of 0.6582. The analysis provides clear, actionable insights for property developers and investors, emphasizing the critical roles of building quality, location, and physical attributes. The conceptual causal inference framework further strengthens these recommendations by highlighting direct levers for business intervention.

### 7.1 Potential Next Steps:

- **Advanced Feature Engineering:** Explore creating more sophisticated features (e.g., distance to nearest metro/mall if coordinates become available, or a "luxury score" for buildings).
- **Hyperparameter Optimization:** Conduct more exhaustive hyperparameter tuning for the XGBoost model using techniques like Grid Search or Randomized Search to potentially eke out even higher performance.
- **Time Series Analysis:** Incorporate the `transaction_date` more deeply to model market trends, seasonality, and predict future price movements.
- **External Data Integration:** Explore the integration of macroeconomic indicators or specific development project data to further enhance predictive power and causal understanding. **SHAP/LIME Implementation:** For individual property pricing decisions, implement SHAP or LIME to provide highly localized and interpretable explanations for each prediction.

## 8 Presentation Notes (for a 15-minute interview)

**Focus:** Clarity, actionable insights, business relevance, and explainability.

### 1. Introduction (1 min):

- Briefly introduce ProjektAnalytics and your role.
- State the chosen sector (Real Estate) and the specific business problem.
- Hook: "Our goal: to help developers price properties optimally and understand what truly drives value in Dubai."

### 2. Data & Key Insights (2 min):

- Mention the dataset source and size.
- Show 1-2 most impactful descriptive plots (e.g., price distribution, price vs. area).
- Highlight average price/area and most common types.
- Briefly mention data cleaning and feature engineering (especially Target Encoding for `building_name`).

### 3. Our Predictive Model (4 min):

- **Model Choice:** "We chose XGBoost with GPU acceleration because it's highly accurate for this type of data, and its GPU capability allowed us to train quickly on a large dataset."
- **Performance:** "The model achieved an R-squared of **65.82%**, meaning it explains over two-thirds of property price variation. Our average prediction error (MAE) is about **1.02 million AED**." **Explainability:** "Crucially, the model tells us *why* it makes predictions. Our feature importance analysis shows that **property size, building name, and area name are the top three drivers of value**." (Show feature importance chart).

### 4. What-If Scenarios (3 min):

- "Let's see how this translates into business decisions."
- Present the base hypothetical property.
- **Scenario 1 (Change Building/Area):** Show the dramatic price change when moving from a premium building/area to a less premium one. "This quantifies the immense value of strategic location and building choice."
- **Scenario 2 (Increase Area):** Show the price change with increased area. Acknowledge the counter-intuitive result if it persists, and frame it as a nuance for luxury properties requiring deeper investigation.
- **Scenario 3 (Increase Rooms):** Show the price change with increased rooms.

### 5. Actionable Insights (4 min):

- "Here are the concrete recommendations for developers and investors:"
- Go through the 4 actionable insights, emphasizing the "Action" part and its direct business impact. Use strong verbs.
- Connect back to the feature importance and what-if scenarios.

### 6. Beyond Correlation: Causal Inference (1 min):

- "To truly understand 'why' things happen, we also considered causal inference."
- Briefly show the conceptual DAG.
- Give one strong example of a **Root Cause, Counterfactual, and Intervention** (e.g., focusing on building quality/brand). "This helps us understand what actions will *cause* a desired outcome, not just what correlates with it."

### 7. Conclusion & Next Steps (Optional, 30 sec):

- "In summary, this model provides powerful, explainable insights for optimizing property pricing in Dubai."
- Briefly mention 1-2 key next steps (e.g., more data, deeper tuning).