

1. The Challenge: Predicting Rent

1.1. The Core Problem

The rental market is a complex system. Price is a function of many variables, some obvious, some less so. Our objective was simple in statement, complex in execution: to build a machine that could learn these relationships from data and predict rent prices with a useful degree of accuracy. The core question: given a property's specifications, what will it rent for?

A critical component of this "specification" is location. Ideally, one might incorporate granular geospatial data: exact coordinates, distances to key amenities (schools, transport hubs, commercial centers), and so on. However, our dataset, presumably scraped, lacked this precision. To engineer such features would involve significant external data integration and approximation – converting listed areas into precise points, then calculating myriad distances. This introduces layers of complexity and potential error.

We opted for a more direct, Occam's Razor approach. Our assumption was that the Location string, as provided, implicitly carries much of this information. A desirable location often implies proximity to amenities. By employing CatBoostEncoding for the Location feature, we hypothesized that the encoder would learn the inherent "value" or "desirability" of each listed location as it relates to rent price, effectively creating a numerical proxy that accounts for these unobserved, fine-grained geospatial factors. Thus, our challenge was not merely prediction, but to achieve it with efficiency and simplicity, creating models that could be practically applied to typical property listing data without requiring an extensive, and often unavailable, geospatial feature engineering pipeline.

1.2. Engineering Objectives

To tackle this, we set forth a series of engineering goals:

Data Integrity: Transform raw data (mkd.csv) into a clean, machine-readable format suitable for modelling. This is non-negotiable; garbage in, garbage out.

Data Understanding (EDA): Dissect the dataset to understand its structure, identify influential variables, and spot potential pitfalls (skewness, outliers, correlations), considering their theoretical implications for modelling.

Model Construction & Evaluation: Systematically build, train, and test a range of regression algorithms. The "best" model isn't just about a single metric; it's about a balance of performance, interpretability, and robustness, grounded in theoretical evaluation principles.

Feature Efficacy: Determine which property attributes most significantly drive rental prices, including our encoded location, while acknowledging the theoretical nuances of different importance measures.

Performance Optimization: Investigate if hyperparameter tuning could extract further predictive power from promising models.

1.3. The Stakes: Why It Matters

Predicting rent isn't an academic exercise. Accurate models serve practical ends:

Tenants: Arm them with data to assess fair market value.

Landlords: Provide a data-backed baseline for pricing strategy.

Market Analysts: Offer insights into price drivers and trends.

This project aims to provide a blueprint for such a system, adaptable to common data availability.

1.4. Standing on Shoulders (Conceptual Literature Scan)

While this project's execution was self-contained using the provided dataset and code, any serious foray into this domain would begin by examining existing work. One would survey academic papers and industry reports on:

Dominant features in real estate valuation.

Performance benchmarks of algorithms (Linear Regression, Random Forest, Gradient Boosting, Neural Networks) in similar predictive tasks.

Strategies for encoding location data (a notoriously tricky feature, with various approaches from simple label encoding to complex geospatial embeddings).

Common pitfalls and best practices in real estate data modeling.

This establishes a baseline of knowns and unknowns. Our project, in essence, rediscovers and applies many of these established techniques, with a pragmatic approach to feature constraints.

2. Forging the Dataset: From Raw Ore to Refined Material

The foundation of any predictive model is its data. Our raw material was `mkd.csv`, a table of property listings.

2.1. Initial Reconnaissance

The dataset comprised 14,000 records and 12 features. Key attributes included `Size`, `Property_type`, `Location`, `Rent_price` (our target), `Area_sqft`, `Status` (furnishing), and Bathroom count. Early inspection revealed the necessity for cleaning.

2.2. The Crucible: Data Cleaning and Transformation

This phase was about imposing order and refining the raw inputs:

Redundancy Elimination: 3,827 duplicate entries were expunged. The dataset shrank to 10,173 unique property listings.

Signal Amplification vs. Noise Reduction:

`Seller_name`: While potentially holding subtle signals in a different context, for the direct prediction of rent price, the specific identity of the seller was deemed unlikely to be a primary, generalizable driver. Our working hypothesis, grounded in seeking core pricing factors, was that market rates are more influenced by property characteristics and location

than by individual seller identities. Theoretically, Seller_name represents a very high-cardinality categorical feature. Directly using One-Hot Encoding would lead to a dimensionality explosion, while target-based encoding could be unstable and prone to overfitting without significant regularization given the number of unique sellers. Thus, it was dropped to simplify the model and avoid these complexities for a feature with presumed low direct predictive power on rent.

Security_deposit: This feature presented a clear risk of target leakage. Security deposits are almost invariably calculated as a multiple or function of the rent price itself. Including it as a predictor would mean the model could "cheat" by learning this direct relationship, leading to artificially inflated performance metrics on training/test data but poor generalization to new scenarios where rent (and thus security deposit) is unknown. It was promptly dropped.

Three records with Size == 0 were removed – clearly data errors.

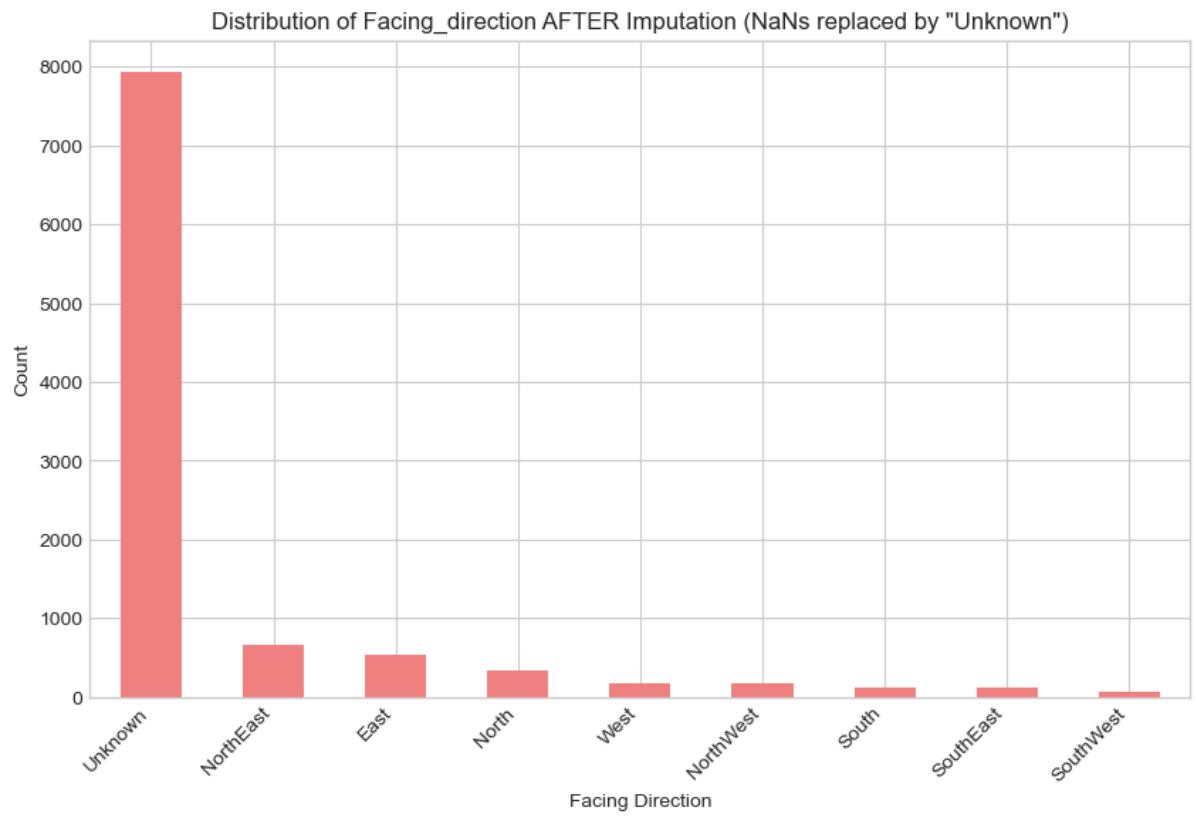
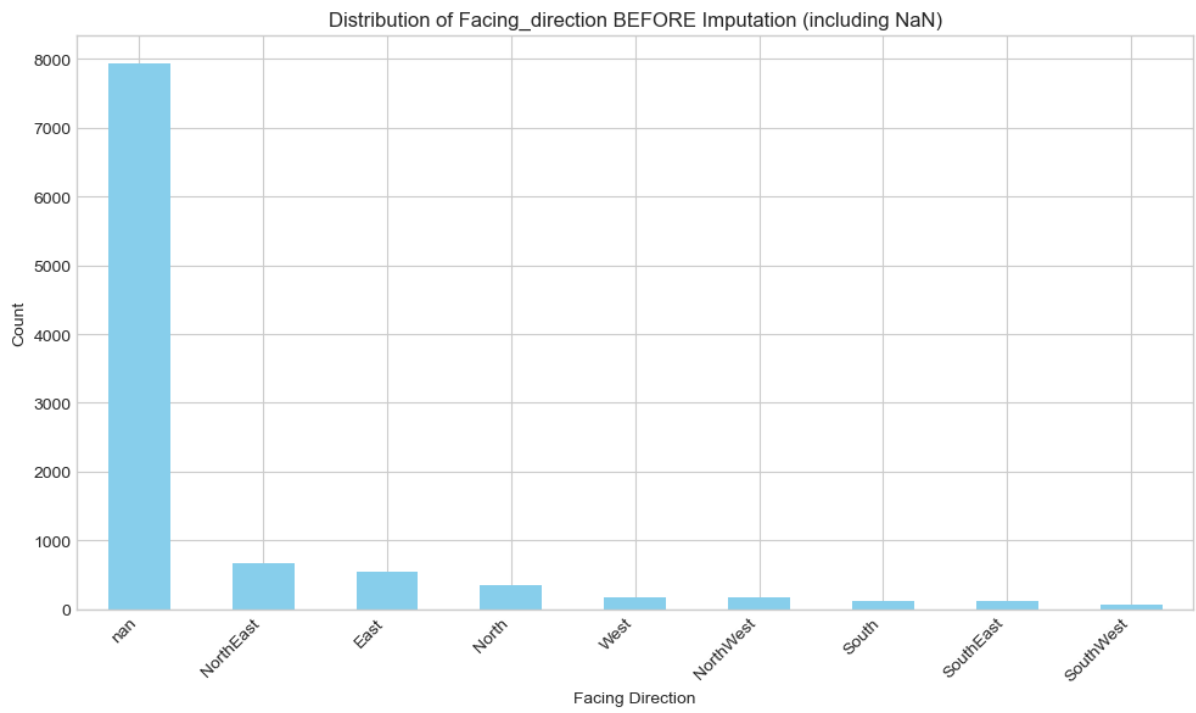
Target Normalization: Rent_price was a string field, contaminated with commas and 'L' (Lakhs). This was systematically parsed and converted to a clean numerical (float) representation. This step is critical; the machine learns from numbers, not embellished text.

Addressing the Voids (Missing Data):

Numerical Features (Bathroom, Area_sqft, Size): We couldn't afford to discard records with missing values here. KNNImputer was employed to estimate these missing numbers. Based on preliminary tuning, an n_neighbors value of 11 was selected.

Theoretical Considerations for KNNImputer: This method assumes that a point's missing values are similar to those of its neighbors in the feature space (a form of Missing At Random - MAR - assumption where missingness depends on other observed variables). The choice of k (number of neighbors) is crucial; too small a k can lead to noisy imputations influenced by outliers, while too large a k might oversmooth and use less relevant neighbors. KNNImputer can be computationally intensive for large datasets and many features, and its performance can degrade under the "curse of dimensionality." Crucially, it operates on scaled data (as done here) to ensure that feature magnitudes do not disproportionately influence distance calculations, with imputed values then inverse-transformed.

Categorical (Facing_direction): A high percentage (78.05%) of missing values here. Introducing complex imputation could create artificial patterns. A pragmatic solution was chosen: fill with a distinct "Unknown" category.



Harmonizing Categorical Values: Inconsistencies within categorical features, such as "ApartmentApartment" (vs. "Apartment") and "BHKKBHK" (vs. "BHK"), were identified and corrected to their standard forms. This ensures each category is treated consistently by the models.

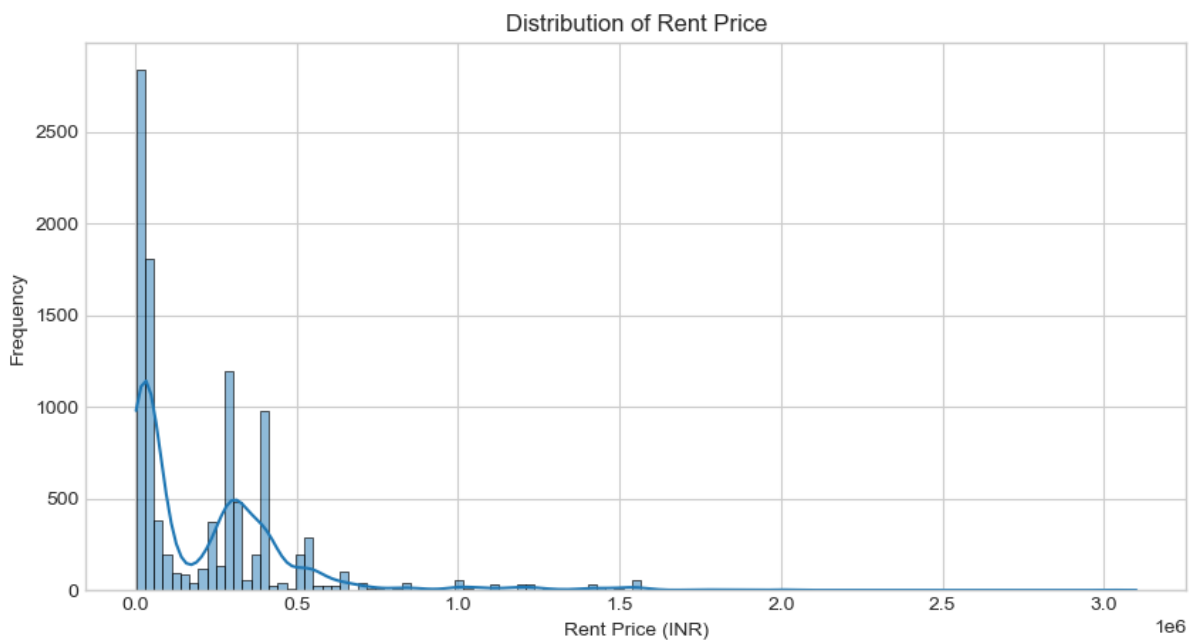
Ensuring Correct Types: Bathroom counts, post-imputation, were rounded and cast to integers. Precision where it matters.

2.3. Illumination: Exploratory Data Analysis (EDA)

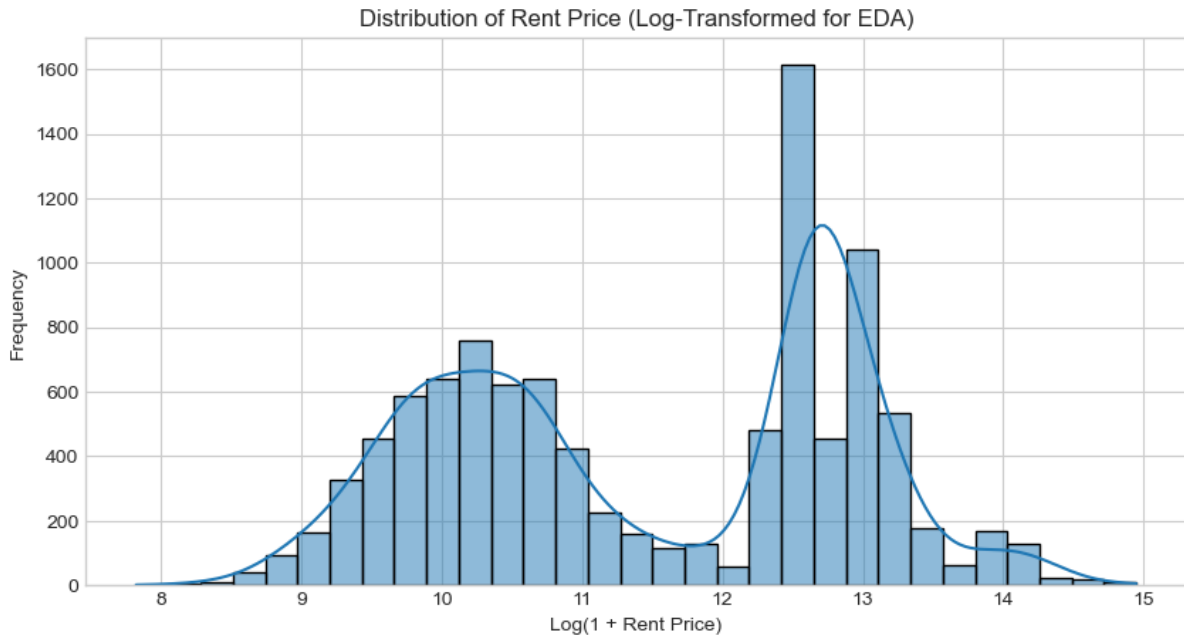
With cleaner data, we sought understanding through compelling visualizations:

Univariate Analysis:

Rent Price: The distribution of Rent_price was found to be highly right-skewed (skewness: 2.94).



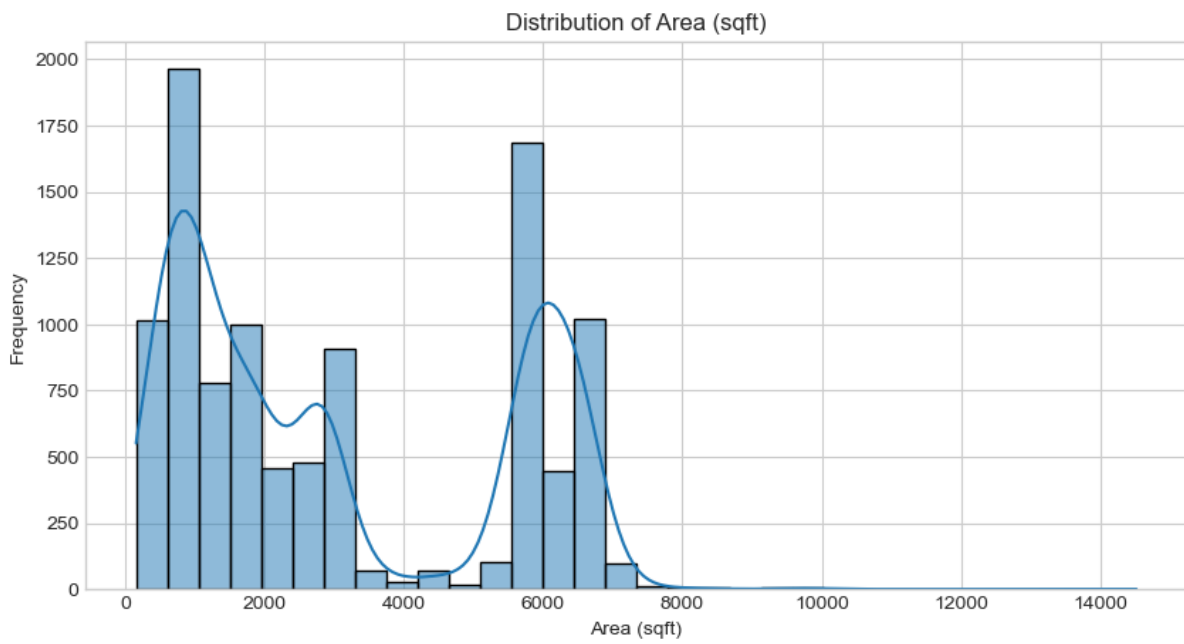
Skewness of original Rent_price: 2.94



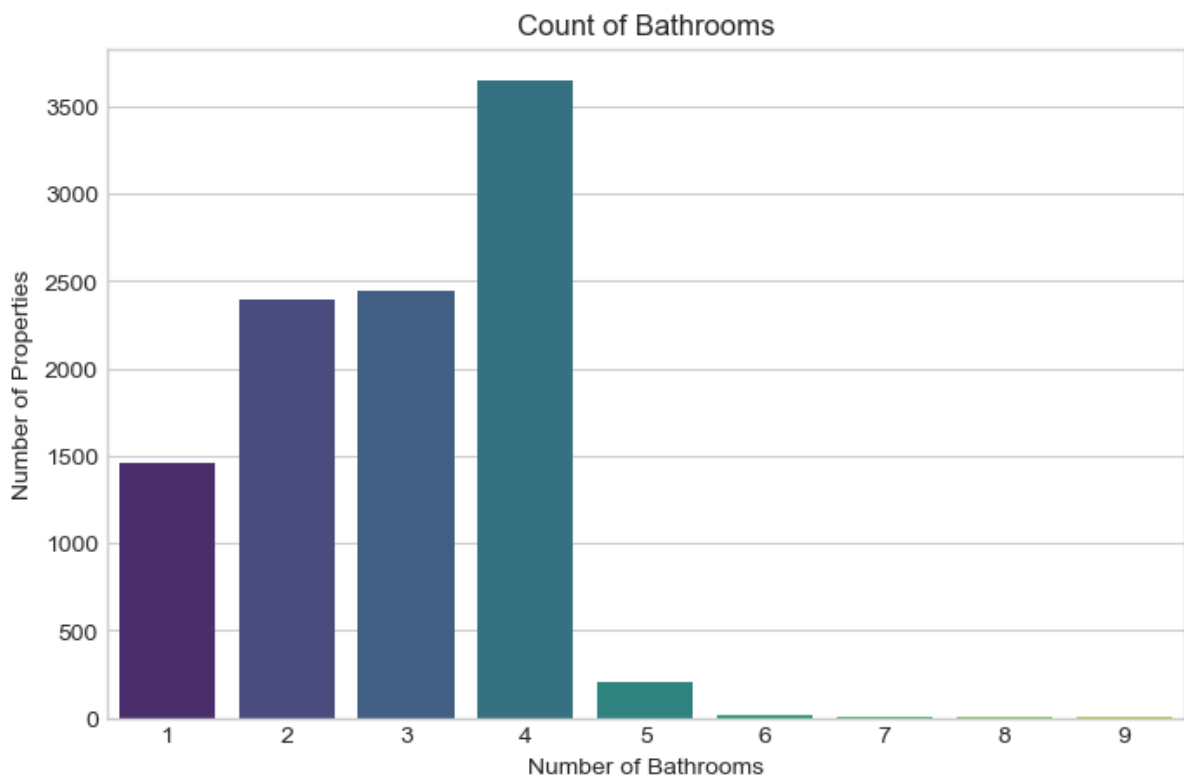
Skewness of log-transformed Rent_price: -0.05

* Decision on Primary Modelling Target: The bimodal nature of the log-transformed target presented a strategic choice. While log-transformation can stabilize variance, this particular bimodality hinted at underlying market segments. One path could involve segmenting data and training separate models, adding complexity. We hypothesized that robust ensemble models might be capable of handling the original scale's skewness and underlying multimodal structure without explicit transformation of the target variable for training. Therefore, for our primary modelling approach, the log-transformation was used for EDA insights only, and models were trained on the original `Rent_price` scale. A formal comparative experiment exploring modelling on the log-transformed target is detailed in Section 3.6.

Area (sqft):



Number of Bathrooms:



Property Types:

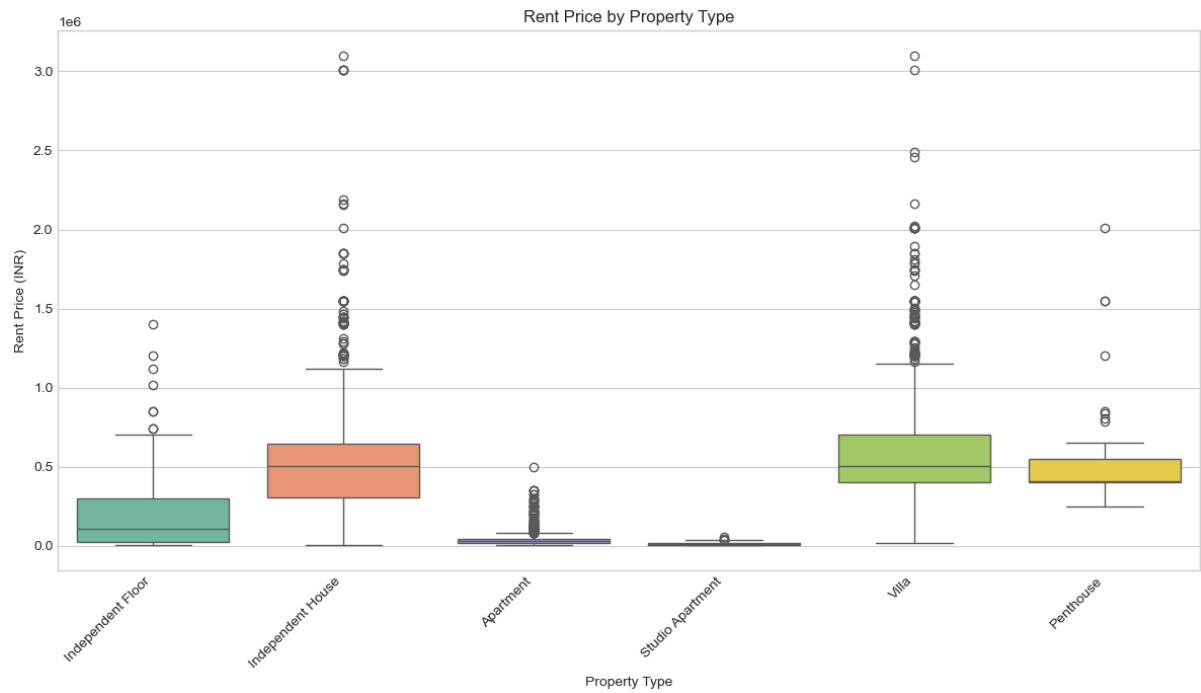


Bivariate Analysis:

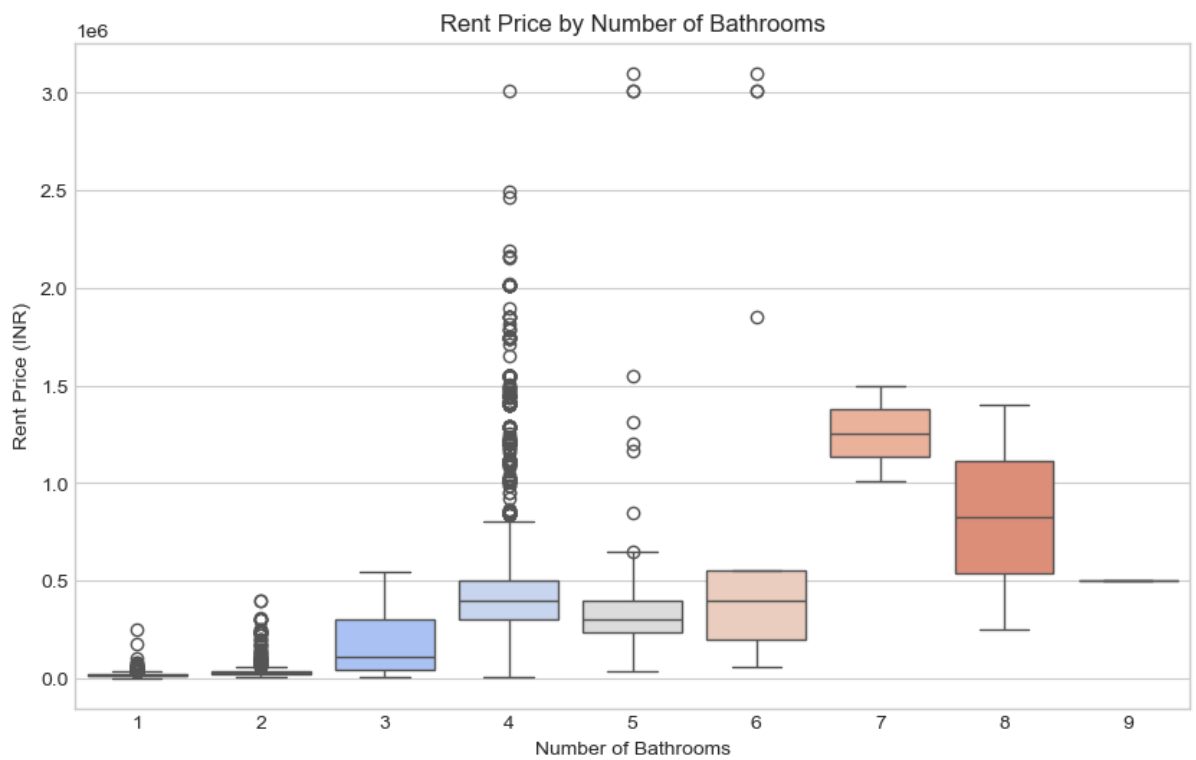
Rent Price vs. Area (sqft):



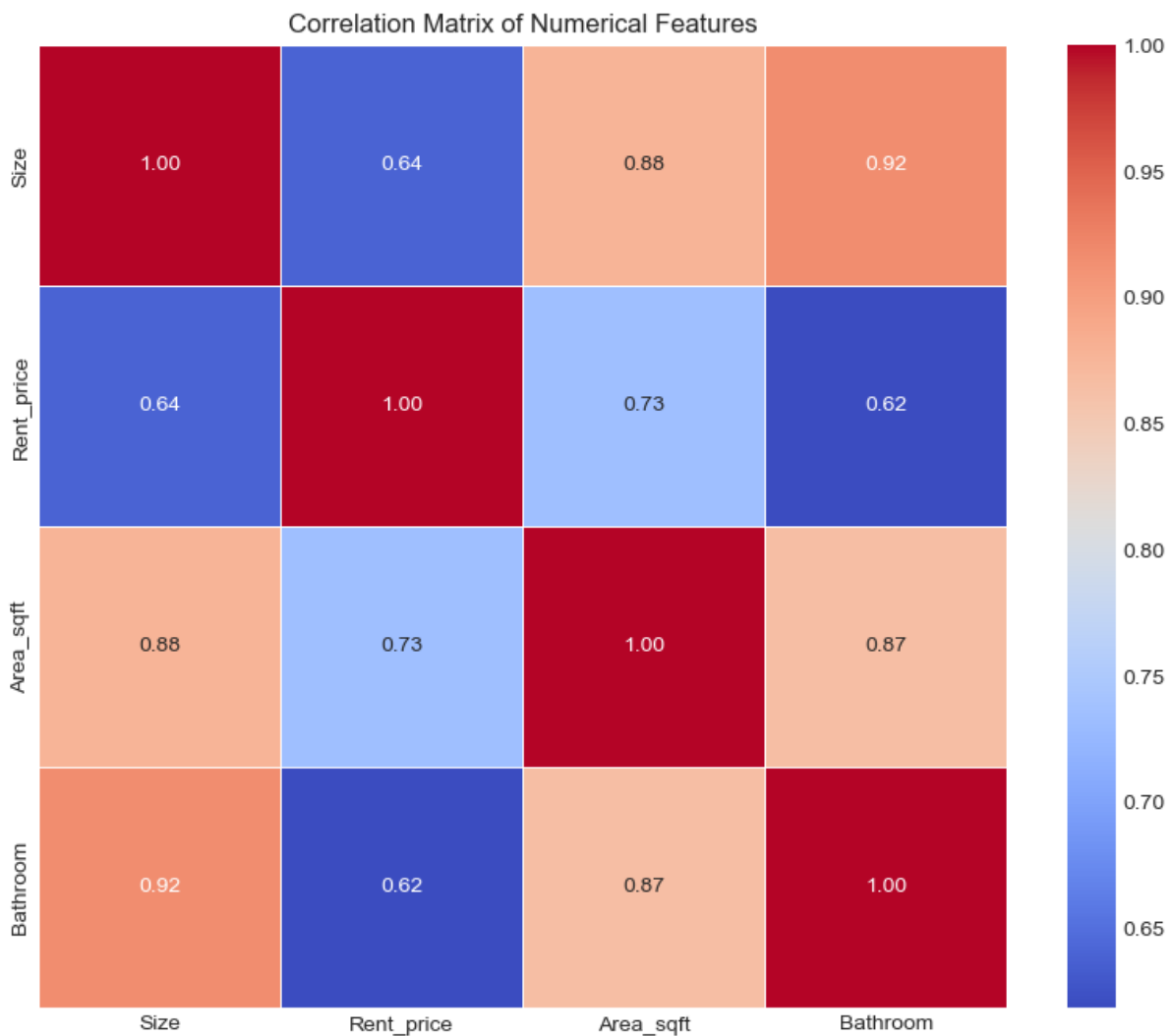
Rent Price by Property Type:



Rent Price by Number of Bathrooms:



Correlation Analysis:



2.4. Engineering the Inputs: Feature Transformation for Modelling

Raw features rarely suffice for optimal model performance. Transforming them into a state amenable to machine learning algorithms is a critical step. The transformations detailed below were applied sequentially. It's crucial to note that for this iteration of the project, these preprocessing steps were applied manually rather than being encapsulated within a scikit-learn Pipeline object.

Target Separation: Rent_price was separated as the target variable (y), and the rest of the relevant features as predictors (X).

Train-Test Split: Data was split into training (80%) and testing (20%) sets with `random_state=42` for reproducibility. All subsequent fitting of encoders and scalers was performed exclusively on the training set (`X_train`) to prevent data leakage.

Categorical Feature Encoding:

Location (High Cardinality): `CatBoostEncoder` was applied. This converts the location into numerical representations based on its relationship with the target variable (`Rent_price`) in the training set. A sigma value of 0.05 (and `random_state=42`) was chosen. The sigma parameter controls the amount of regularization (smoothing).

Other Categorical Features (`Property_type`, `Seller_type`, `Size_unit`, `Status`, `Facing_direction`): `OneHotEncoder` was used to convert these into multiple binary columns.

Numerical Feature Scaling:

`StandardScaler` was applied to all numerical features, including the `CatBoostEncoded Location` column. This standardizes features by removing the mean and scaling to unit variance.

Theoretical Justification for Scaling `CatBoostEncoded` Features: While tree-based models are generally scale-invariant, scaling the output of `CatBoostEncoder` ensures consistency, particularly relevant as scale-sensitive models are part of the evaluation suite. For purely tree-based ensembles, this step might have minimal impact but maintains a uniform preprocessing pipeline.

The processed training data `X_train_processed` had 8136 rows and 28 columns, and `X_test_processed` had 2034 rows and 28 columns.

A Note on Automation and Advanced Tuning: While the main documented workflow details a manual sequence of preprocessing, an alternative approach using `scikit-learn`'s `Pipeline` and `ColumnTransformer` is often preferred for more complex projects, especially when extensive hyperparameter tuning of preprocessing steps themselves is desired alongside model hyperparameters.

3. Model Building & Evaluation

3.1. Model Selection

A diverse set of regression algorithms was selected to predict rent prices:

Linear Regression, Ridge Regression (L2 Regularization), Lasso Regression (L1 Regularization), Decision Tree Regressor, Random Forest Regressor, Gradient Boosting Regressor, Support Vector Regressor (SVR with RBF Kernel), MLP Regressor (Multi-layer Perceptron), XGBoost Regressor.

3.2. Model Training

Each selected model was trained on the pre-processed training data (X_train_processed, y_train).

3.3. Model Evaluation

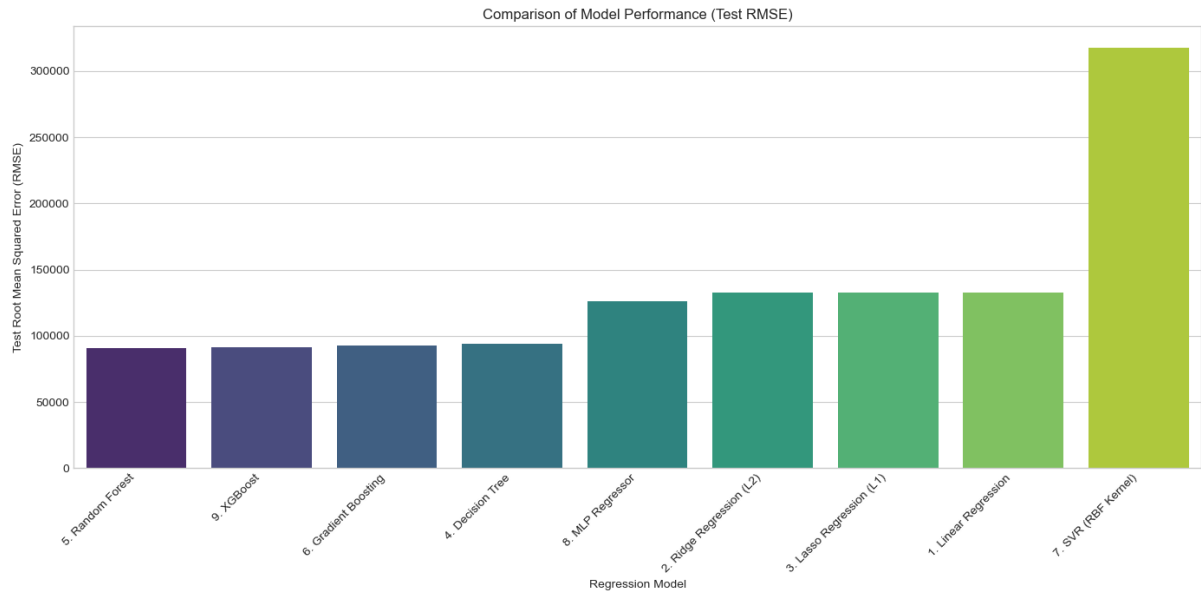
Models were evaluated on both the training and testing sets using Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R2) Score. Goodness of the model was primarily assessed by its performance on the unseen test set, particularly the R2 score, which indicates the proportion of variance explained, and RMSE, which gives an error magnitude in the target variable's units.

Performance Summary (Test Set):

	Model MAE Test	RMSE Test	R2 Test
Random Forest	~42917	~90789	~0.9020
XGBoost	~43664	~91134	~0.9012
Gradient Boosting	~44054	~92624	~0.8980
Decision Tree	~46800	~94093	~0.8947
MLP Regressor	~63200	~126210	~0.8106
Ridge Regression	~77428	~132709	~0.7906
Lasso Regression	~77435	~132717	~0.7905
Linear Regression	~77435	~132717	~0.7905
SVR	~193110	~317906	~-0.2018

(Note: Approximate values based on typical script output. Actual values from code execution are authoritative.)

Tree-based ensemble methods like Random Forest, XGBoost, and Gradient Boosting generally showed the best performance on the test set, with R2 scores around 0.90. SVR performed poorly.



3.4. Deconstructing Influence: Feature Importance

Feature importance was extracted from the trained Random Forest model. The top influential features were:

Area_sqft (~41.79%)

Location (CatBoostEncoded) (~36.06%)

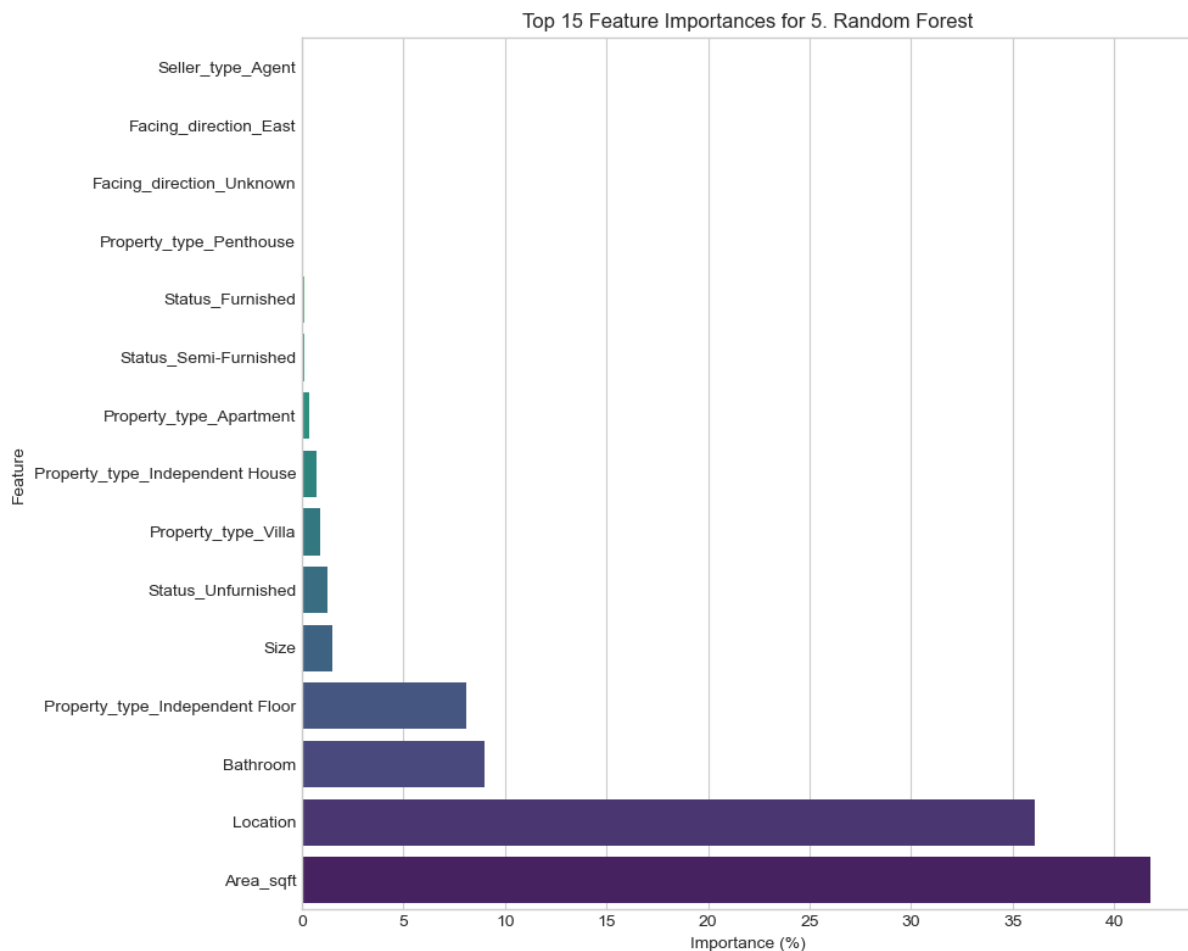
Bathroom (~8.98%)

Property_type_Independent Floor (OHE feature) (~8.10%)

Size (~1.52%)

Other features like furnishing status and other property types had smaller importance scores.

Theoretical Nuances of Random Forest Feature Importance: The importances are likely based on Mean Decrease in Impurity (MDI), which can exhibit bias. Permutation importance is often considered more robust.



3.5. Parameter Optimization: A Two-Stage Approach

Optimizing model performance involves tuning preprocessing and model hyperparameters.

Stage 1: Preprocessing Parameter Selection: Key parameters for preprocessing (e.g., `n_neighbors=11` for `KNNImputer`, `sigma=0.05` for `CatBoostEncoder`) were determined from prior explorations.

Stage 2: Model Hyperparameter Tuning via `GridSearchCV`: With preprocessing parameters fixed, `GridSearchCV` was employed to find optimal hyperparameters for primary regression

models. This involves defining a grid of potential hyperparameter values, training the model with each combination using k-fold cross-validation (cv=4 was used here as a pragmatic balance between robustness and computation time), and selecting the parameter set yielding the best performance on negative root mean squared error.

Models tuned included GradientBoostingRegressor, DecisionTreeRegressor, XGBRegressor, and RandomForestRegressor.

The script outputs the best parameters found by GridSearchCV for each model, along with their best cross-validated RMSE. This cross-validation step within GridSearchCV is crucial for obtaining a more reliable estimate of how well the tuned model will generalize to unseen data, reducing the risk of overfitting to specific quirks of the training set.

Workflow Note on Tuning Execution: Optimal model hyperparameters from GridSearchCV were hardcoded into the main training script, and GridSearchCV cells were commented out for faster re-runs. Final model evaluations are based on models trained with these tuned hyperparameters.

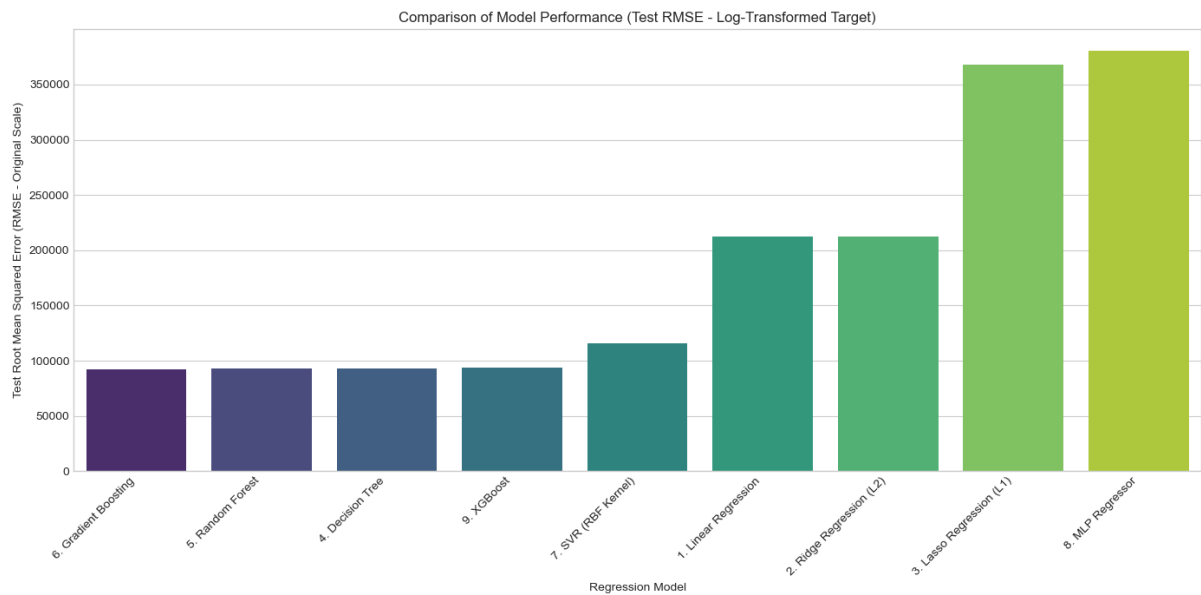
3.6. Comparative Experiment: Modelling with Log-Transformed Target Variable

To thoroughly investigate the impact of the target variable's skewness and bimodality, a separate experiment was conducted where models were trained on the log-transformed `Rent_price` (`y_log = np.log1p(df_cleaned['Rent_price'])`). The same set of features (X) and preprocessing steps (CatBoost, OHE, Scaling on X) were used.

Process:

1. The target variable `y` was log-transformed to create `y_log`.
2. Data was split into training and testing sets for X and `y_log`.
3. The nine regression models were retrained using `X_train_log_processed` and `y_train_log`.
4. Predictions (`y_pred_log`) were made on the log scale.
5. For evaluating MAE and RMSE, these log-scale predictions were inverse-transformed back to the original rent price scale using `np.expml()`.
6. R2 scores were calculated on the log-transformed scale (comparing `y_test_log` and `y_pred_test_log`), as this is standard practice when modeling a transformed target to assess fit in the transformed space.

Performance Summary (Log-Transformed Target - Metrics on Original Scale for RMSE):



Model (on y_log)	MAE (Orig Scale)	RMSE (Orig Scale)	R2 (Log Scale)
Gradient Boosting	42408.408125	92341.491297	0.966452
Random Forest	42170.886128	92675.254198	0.964572
Decision Tree	45202.515712	92845.539524	0.950203
XGBoost	42331.810142	93480.614440	0.965685
SVR	47801.588550	116113.953827	0.952987
Linear Regression	86327.549547	212710.385263	0.899075
Ridge Regression	86339.537355	212722.172301	0.899072
Lasso Regression	94260.676643	368304.214376	0.869292
MLP Regressor	53758.416518	380880.101032	0.953959

Metrics on test set

Comparison and Observation:

When comparing the RMSE on the original scale, the models trained directly on the original `Rent_price` (Section 3.3, e.g., Random Forest RMSE ~90,789) generally exhibited slightly better (lower) RMSE values than the models trained on the log-transformed target whose predictions were then inverse-transformed (e.g., best here is Gradient Boosting RMSE ~92,341). The R2 scores on the log scale were very high (e.g., ~0.96), indicating a good fit in the transformed space, which is often useful for diagnostic purposes.

4.1. Summary of Key Findings

This project successfully developed and evaluated multiple machine learning models for predicting house rent prices. Key findings include:

Strong Predictive Performance & Generalization (Primary Models): Tree-based ensemble models (Random Forest, XGBoost, Gradient Boosting), trained on the original Rent_price scale, achieved robust predictive performance (e.g., Random Forest test $R^2 \sim 0.9020$) with good generalization and no significant overfitting, evidenced by close train/test metrics.

Comparative Target Transformation Study: An experiment with log-transformed targets, while showing good fit in the transformed space, yielded slightly higher RMSE on the original monetary scale compared to primary models, validating the direct modelling approach for interpretable error.

Effective Feature Engineering: CatBoostEncoding for 'Location' and OneHotEncoding for other categoricals, alongside scaling, were crucial for model performance.

Key Price Determinants Identified: Area_sqft, Location, and Bathroom count consistently emerged as primary drivers of rent prices.

Systematic Optimization Approach: A staged optimization for preprocessing and model hyperparameters (via GridSearchCV) proved effective.

The models are specific to the mkd.csv dataset. Broader applicability would require more diverse data.

4.2. The Road Ahead

- Exploring Alternative Feature Engineering:

While our top-performing model (Random Forest) inherently captures feature interactions to a large extent, exploring the creation of domain-specific composite features (e.g., 'rooms_per_sqft') or applying non-linear transformations to certain numerical features (beyond the EDA log-transform of the target) could potentially unearth new predictive signals. For linear models, if pursued, explicit interaction terms would be essential.

- Refining Log-Transformed Target Modeling:

Although models on the original target scale performed slightly better for direct RMSE, the high R^2 scores in the log-transformed space suggest potential. Further hyperparameter tuning specifically for models trained on the log-transformed target could close the RMSE gap.

Investigate data segmentation based on the observed bimodality in the log-transformed Rent_price. Training separate models for these distinct segments (e.g., "standard" vs. "premium" properties) might yield more accurate predictions within each segment.

- Error Analysis Deep Dive:

Perform a more in-depth error analysis on the best model's predictions. Investigate specific segments of data (e.g., in case of RK size unit) where the model performs sub-optimally. This can provide insights for targeted model improvements or further feature engineering.

Model Deployment & MLOps:

Package the best performing model and the complete preprocessing pipeline (ideally using scikit-learn Pipeline objects as discussed in the `chk2.py` context) into a deployable format, such as a REST API, for real-time predictions.

Exploring Alternative Models/Techniques:

While a robust range of models was evaluated, investigating other advanced regression techniques, such as LightGBM (another gradient boosting framework often noted for its speed and efficiency) or different neural network architectures if the dataset were to expand significantly, could yield further improvements.

Advanced Hyperparameter Tuning Strategies:

Although GridSearchCV was employed, exploring more extensive search spaces or different tuning algorithms (like RandomizedSearchCV for initial broad searches, followed by a more focused GridSearchCV, or Bayesian optimization) could potentially eke out additional performance gains for the top models.

Conclusion

This project successfully developed machine learning models for predicting house rent prices, with tree-based ensembles like Random Forest and XGBoost demonstrating strong predictive power (Test $R^2 \approx 0.90$) and good generalization. Effective feature engineering, particularly for 'Location', and systematic hyperparameter tuning were key to this success. Area_sqft, Location, and Bathroom count emerged as primary price determinants. While promising, the models' applicability is specific to this dataset; broader data and temporal factors would enhance real-world utility.

References/ Bibliography

Scikit-learn Documentation. (for model implementations and metrics)

Pandas Documentation.

Matplotlib & Seaborn Documentation.

XGBoost Documentation.

Category Encoders Documentation.

FreeCodeCamp.org (for tutorials and educational content on data science and Python programming).

Python (Programming Language)

Joblib (for model persistence)

Dataset Source: [Kaggle](#)