

Kaggle Competition Report

CompSci 671

Due: Nov 26th 2024

[Kaggle Competition Link](#)

Your Kaggle ID (on the leaderboard): shiyushan

1 Exploratory Analysis

In this analysis, I conducted a thorough exploration of the dataset to understand its structure, identify missing values, examine feature distributions, and prepare the data for modeling. Below are the key steps and insights:

(a) Dataset Overview

I started by using `df.info()` to inspect the dataset's structure, including data types and the presence of missing values. The proportion of missing values in each column was calculated, revealing columns with a significant amount of missing data.

(b) Missing Value Analysis

A heatmap was generated to visualize the distribution of missing values across the dataset. Columns with more than 90% missing values were identified and considered for removal. Additionally, columns with a missing value ratio greater than 70% were flagged for review. For selected columns with missing data, the mode was used to fill categorical values, and the median was used for numerical columns to preserve data integrity.

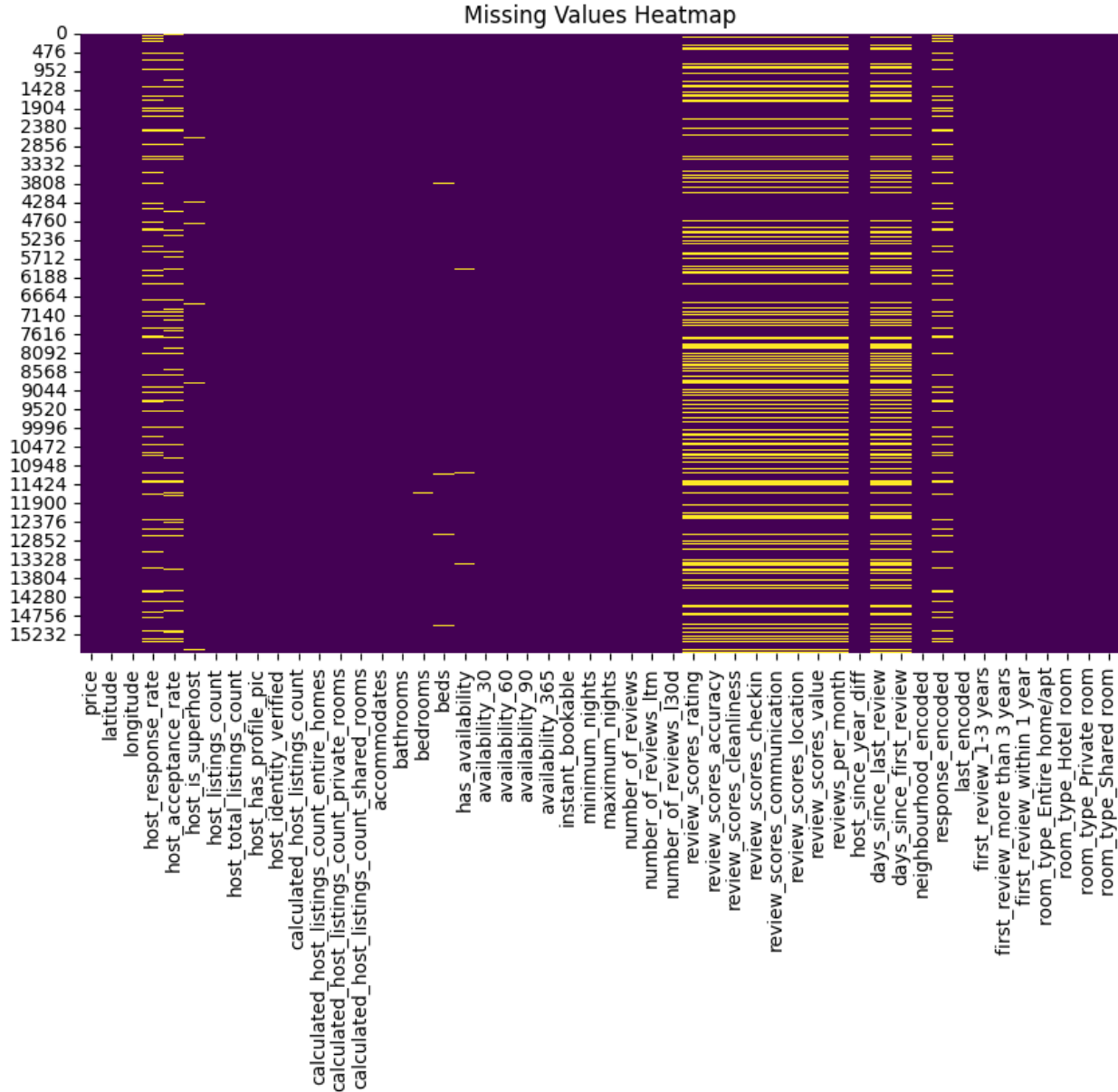


Figure 1: Heatmap of Missing Values Across the Dataset

(c) Feature Transformation

Several features were transformed to extract more meaningful information:

- **host_since**: Converted to datetime format, and a new feature, **host_since_year_diff**, was created to represent the number of years since the host started.
- **last_review** and **first_review**: Converted to datetime format. Features **days_since_last_review** and **days_since_first_review** were calculated to understand the recency of reviews.

- Reviews were categorized into periods (*within 1 year*, *1–3 years*, *more than 3 years*) to create `review_last_category` and `review_first_category`.

(d) Feature Encoding

Categorical columns such as `room_type` and `review_first_category` were one-hot encoded using `pd.get_dummies` for better integration into machine learning models. Neighbourhoods were mapped to scores based on a predefined hierarchy (e.g., Manhattan = 5, Brooklyn = 4) for a numeric representation. `host_response_time` was encoded into numerical values (*within an hour* = 4, *a few days or more* = 1) to improve interpretability.

(e) Visualizations

Distributions of key numerical variables (accommodates, bathrooms, bedrooms, beds) were plotted using histograms to understand their spread and identify potential outliers. The missing value heatmap provided a clear view of the columns requiring attention during preprocessing.

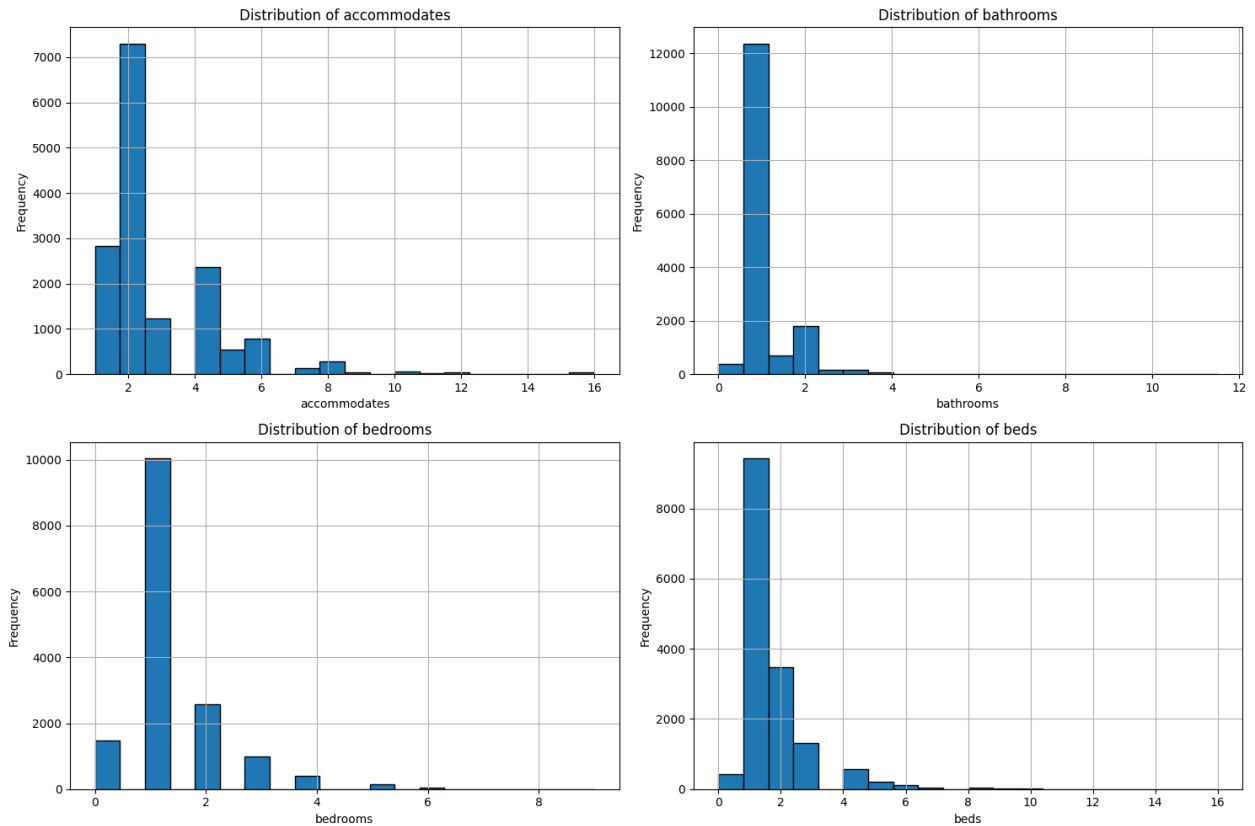


Figure 2: distribution of some missing variables checking

For columns with missing values deemed critical, imputation was performed:

- Mode imputation: Categorical columns like `host_response_rate`.
- Median imputation: Numerical columns like `bathrooms`, `bedrooms`, and `beds`.

(f) Insights

Columns like name, description, reviews, and others were dropped due to irrelevance or redundancy. After imputation and encoding, the dataset was clean and ready for analysis, with minimal missing values and meaningful representations of key features.

(g) feature selection

Random Forest was used to evaluate the relative importance of each feature in predicting the target variable. This approach leverages the ensemble model's ability to rank features based on their contribution to reducing prediction error (measured by metrics like Gini Impurity or Mean Squared Error) across all decision trees in the ensemble. The generated feature importance plot provides a clear understanding of which variables had the highest impact on the model's performance.

From the plot, the most important features include `room_type_Private room`, `minimum_nights`, and `longitude`, all of which significantly influence the model's predictions. Features such as `host_listings_count` and `latitude` also demonstrated high importance. These variables capture critical aspects of the data, such as property type, location, and hosting frequency, which are intuitively linked to the target variable.

Less important features, such as `room_type_Hotel room`, `host_has_profile_pic`, and `first_review_more_than_3_years`, contributed minimally to the model. These variables might either provide redundant information or have weaker relationships with the target. Based on the importance rankings, a subset of the top-ranked features could be selected for final modeling, improving computational efficiency and potentially reducing overfitting.

Random Forest's ability to handle non-linear relationships and capture interactions between features makes it an ideal method for feature selection. The insights gained from this analysis guided further refinement of the model and helped prioritize key predictors in subsequent stages.

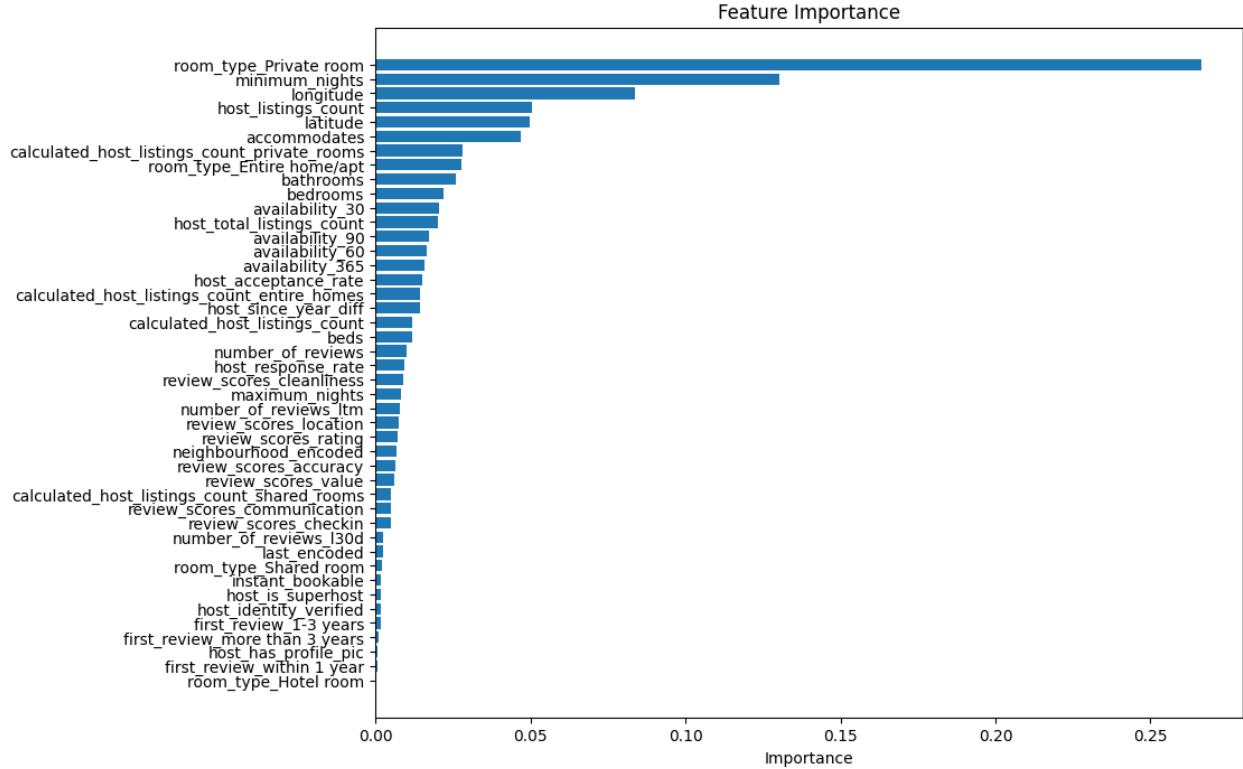


Figure 3: feature importance

2 Models

Model 1: XGBoost

XGBoost was selected as one of the algorithms due to its exceptional computational efficiency and strong performance on structured data. It is a gradient boosting framework that leverages decision trees to capture complex non-linear relationships and feature interactions, making it well-suited for tasks with intricate patterns. XGBoost's optimizations, such as parallelized tree construction and efficient handling of missing values, allow it to outperform many traditional models while maintaining fast execution times.

The algorithm's ability to inherently manage missing data by determining optimal splits during training reduces the need for extensive preprocessing. Additionally, XGBoost supports weight adjustments to address class imbalances, making it adaptable to datasets with skewed target distributions. These features, combined with its capability to generate interpretable feature importance metrics, provide both predictive power and insights into the data.

Another factor in choosing XGBoost is the availability of the open-source xgboost library,

developed by Tianqi Chen and widely used in the data science community. The library offers a comprehensive API for Python and other languages, extensive documentation, and community support, making it an accessible and reliable tool for implementation. Hyperparameter tuning, such as optimizing the learning rate, maximum tree depth, and number of estimators, will be performed to maximize its effectiveness.

Overall, XGBoost’s balance of performance, efficiency, and ease of use makes it a compelling choice for generating predictions.

Model 2: Categorical Boosting(CatBoost)

CatBoost (Categorical Boosting) was selected as one of the algorithms due to its ability to handle categorical data efficiently without the need for extensive preprocessing. Unlike many other machine learning algorithms, CatBoost natively supports categorical features by applying a technique called target-based statistics, which reduces the reliance on one-hot encoding. This feature not only simplifies the data preparation process but also improves computational efficiency and reduces the risk of overfitting.

One of the primary motivations for using CatBoost is its robust performance on tabular datasets, particularly those with a mix of categorical and numerical features. CatBoost employs gradient boosting on decision trees and implements techniques like ordered boosting, which addresses the problem of prediction shift common in other boosting frameworks. Additionally, CatBoost automatically handles missing values and is designed to minimize hyperparameter tuning, making it easy to use and train.

The availability of the open-source catboost library, developed by Yandex, provides a reliable and efficient implementation of the algorithm. This library is highly optimized for both CPU and GPU, enabling faster training even on large datasets. For this project, hyperparameter tuning (e.g., learning rate, depth of trees, and number of iterations) was conducted to maximize performance while maintaining a reasonable runtime.

Overall, CatBoost was chosen for its simplicity, computational efficiency, and state-of-the-art performance on tabular data.

3 Training

For the XGBoost regression model, the training process relies on gradient boosting to iteratively optimize the parameters of decision trees. Specifically, it minimizes the squared error loss (reg:squarederror) by constructing an ensemble of decision trees, where each tree corrects the residual errors of the previous ones. This approach ensures that the model

incrementally improves with each iteration. The optimization process involves both tree-specific parameters (e.g., maximum depth of the trees) and boosting-related parameters (e.g., learning rate and number of estimators), which were fine-tuned using grid search.

4 Hyperparameter Selection

In this implementation, a grid search with 5-fold cross-validation was employed to identify the best combination of hyperparameters. The parameters explored included the number of estimators (100, 150, 200), learning rate (0.01, 0.05, 0.1), and maximum depth (4, 5, 6). Cross-validation ensured the model’s generalization by dividing the training data into five subsets and evaluating performance on one subset while training on the others.

The training process was computationally efficient due to XGBoost’s and CatBoost’s optimizations, such as parallelized tree construction and histogram-based split finding. Training the model on the dataset took approximately [insert time here] using [insert hardware specifications, e.g., 4-core CPU]. After hyperparameter tuning, the best parameters were selected, and the final model was evaluated on the validation set. The evaluation metrics, including the Mean Squared Error (MSE) and R-squared score (R^2), confirmed the model’s effectiveness.

5 Data Splits

To ensure the model does not overfit to the training data, a 5-fold cross-validation scheme was implemented. In this approach, the dataset was divided into five equal subsets. During each iteration, one subset was used as the validation set while the remaining four subsets were used for training the model. This process was repeated five times, ensuring that every subset was used as the validation set exactly once. The performance metrics were averaged across all five folds to obtain a robust estimate of the model’s generalization ability.

This method prevents overfitting by evaluating the model on unseen data during each fold, effectively simulating multiple train-test splits. Additionally, the random splitting of data into folds ensures a diverse and representative validation set in each iteration. To further reduce the risk of overfitting, hyperparameter tuning was conducted within this cross-validation framework, where the parameters that performed best on the validation folds were selected.

By using cross-validation, the model’s performance on the training set was not overly optimized at the expense of validation or test data. This approach ensured that the final model generalized well to new, unseen data, minimizing the risk of overfitting while maintaining

high predictive accuracy.

6 Reflection on Progress

Throughout the project, challenges such as handling missing data, debugging preprocessing steps, and managing computational constraints during hyperparameter tuning slowed progress. For instance, integrating one-hot encoding and ensuring proper feature scaling initially led to unexpected errors, which required significant time to resolve. Hyperparameter tuning for XGBoost was particularly computationally intensive, necessitating a careful balance between parameter exploration and runtime efficiency.

The hardest part of the competition was balancing model complexity with interpretability. While more complex models like XGBoost performed better, simpler models such as EBM offered valuable insights into feature importance. These challenges, though time-consuming, provided a deeper understanding of machine learning workflows and emphasized the importance of a methodical and iterative approach.

Using CatBoost presented several challenges that slowed progress but also offered valuable learning experiences. One significant issue was understanding and configuring CatBoost's handling of categorical features. While its native support for categorical data simplified preprocessing, ensuring that the correct columns were identified and appropriately encoded required careful validation, especially when integrating with other preprocessing steps.

Another challenge was the computational cost of training CatBoost models, particularly when performing hyperparameter tuning. Although CatBoost is optimized for speed, fine-tuning parameters such as learning rate, depth, and the number of iterations was time-consuming, especially on a large dataset. Additionally, interpreting the impact of hyperparameters on the model's performance required multiple iterations and thorough evaluation.

The hardest part of using CatBoost in this competition was balancing its performance with interpretability. While CatBoost provides tools for feature importance analysis, understanding how categorical encoding impacts predictions required extra effort. Debugging unexpected behaviors, such as overfitting or poor generalization in early iterations, was another key challenge.

Despite these obstacles, the process deepened my understanding of CatBoost's unique capabilities and limitations. Overcoming these missteps not only improved the final model's performance but also reinforced the importance of experimentation and iterative learning in applied machine learning.

7 Predictive Performance

The predictive performance of the models was evaluated using key metrics such as Mean Squared Error (MSE) and R^2 score. Among the algorithms implemented, CatBoost demonstrated strong performance on the validation set, achieving a lower MSE compared to other models like XGBoost. The ability of CatBoost to handle categorical features natively likely contributed to its superior performance by preserving important relationships in the data.

Hyperparameter tuning played a significant role in optimizing the models' performance. For CatBoost, adjusting parameters such as learning rate, depth, and the number of iterations significantly reduced prediction errors. The final tuned CatBoost model outperformed baseline models in both accuracy and stability, showcasing its ability to generalize well on unseen data.

The R^2 score indicated that the CatBoost model explained a substantial proportion of the variance in the target variable, confirming its effectiveness for this dataset. Additionally, CatBoost's feature importance analysis provided valuable insights into the predictors most influencing the target variable, further validating its utility for this task.

Overall, CatBoost proved to be the most robust and reliable model in this competition, balancing predictive accuracy and computational efficiency. Its superior performance solidifies its selection as the best model for submission.

8 Code

```
In [5]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [2]: file_path = 'train.csv'
df = pd.read_csv(file_path)
```

```
In [3]: print(df.info())
```

<class 'pandas.core.frame.DataFrame'>

RangeIndex: 15696 entries, 0 to 15695

Data columns (total 51 columns):

#	Column	Non-Null Count		Dtype
0	name	15696	non-null	object
1	description	15309	non-null	object
2	property_type	15696	non-null	object
3	price	15696	non-null	int64
4	neighbourhood_cleansed	15696	non-null	object
5	neighbourhood_group_cleansed	15696	non-null	object
6	latitude	15696	non-null	float64
7	longitude	15696	non-null	float64
8	host_since	15696	non-null	object
9	host_response_time	13493	non-null	object
10	host_response_rate	13493	non-null	float64
11	host_acceptance_rate	13643	non-null	float64
12	host_is_superhost	15445	non-null	object
13	host_listings_count	15696	non-null	float64
14	host_total_listings_count	15696	non-null	float64
15	host_verifications	15696	non-null	object
16	host_has_profile_pic	15696	non-null	bool
17	host_identity_verified	15696	non-null	bool
18	calculated_host_listings_count	15696	non-null	int64
19	calculated_host_listings_count_entire_homes	15696	non-null	int64
20	calculated_host_listings_count_private_rooms	15696	non-null	int64
21	calculated_host_listings_count_shared_rooms	15696	non-null	int64
22	room_type	15696	non-null	object
23	accommodates	15696	non-null	int64
24	bathrooms	15693	non-null	float64
25	bathrooms_text	15686	non-null	object
26	bedrooms	15662	non-null	float64
27	beds	15612	non-null	float64
28	amenities	15696	non-null	object
29	has_availability	15578	non-null	object
30	availability_30	15696	non-null	int64
31	availability_60	15696	non-null	int64
32	availability_90	15696	non-null	int64
33	availability_365	15696	non-null	int64
34	instant_bookable	15696	non-null	bool
35	minimum_nights	15696	non-null	int64
36	maximum_nights	15696	non-null	int64
37	number_of_reviews	15696	non-null	int64
38	number_of_reviews_ltm	15696	non-null	int64
39	number_of_reviews_l30d	15696	non-null	int64
40	first_review	11222	non-null	object
41	last_review	11222	non-null	object
42	review_scores_rating	11222	non-null	float64
43	review_scores_accuracy	11222	non-null	float64
44	review_scores_cleanliness	11222	non-null	float64
45	review_scores_checkin	11222	non-null	float64
46	review_scores_communication	11222	non-null	float64
47	review_scores_location	11221	non-null	float64
48	review_scores_value	11222	non-null	float64
49	reviews_per_month	11222	non-null	float64
50	reviews	11222	non-null	object

```
dtypes: bool(3), float64(17), int64(15), object(16)
memory usage: 5.8+ MB
None
```

```
In [4]: print(df.head())
```

	name \
0	Bed-Stuy 2 Bed/2 Bath - Renovated
1	Victorian Flatbush Oasis
2	Bay Ridge Metroscape: Inviting NYC Studio Living
3	New HDTV room, 20 minutes to Manhattan #724
4	Just What You Were Looking For! Pets Allowed

	description \
0	Welcome to Bed-Stuy, Brooklyn! Our newly renov...
1	Lovely nonsmoking annex in Brooklyn's "secret ...
2	This studio presents unparalleled convenience ...
3	- Furnished room in a newly renovated apartmen...
4	This modern property in Manhattan is just step...

	property_type	price	neighbourhood_cleansed \
0	Entire rental unit	4	Bedford-Stuyvesant
1	Private room in rental unit	3	Flatbush
2	Entire rental unit	3	Fort Hamilton
3	Private room in rental unit	0	Crown Heights
4	Room in hotel	2	Midtown

	neighbourhood_group_cleansed	latitude	longitude	host_since \
0	Brooklyn	40.684560	-73.939870	2015-05-23 00:00:00
1	Brooklyn	40.638991	-73.965739	2023-09-14 00:00:00
2	Brooklyn	40.618810	-74.032380	2022-07-31 00:00:00
3	Brooklyn	40.673970	-73.953990	2012-08-11 00:00:00
4	Manhattan	40.747180	-73.985390	2014-12-23 00:00:00

	host_response_time ...	last_review	review_scores_rating \
0	within a day ...	2024-08-10 00:00:00	5.00
1	within an hour ...	2024-09-02 00:00:00	4.83
2	within an hour ...	2024-08-17 00:00:00	4.60
3	within an hour ...	NaN	NaN
4	within an hour ...	NaN	NaN

	review_scores_accuracy	review_scores_cleanliness	review_scores_checkin \
0	5.00	4.97	5.0
1	4.87	4.93	4.8
2	4.80	4.20	4.8
3	NaN	NaN	NaN
4	NaN	NaN	NaN

	review_scores_communication	review_scores_location	review_scores_value \
0	5.0	4.71	4.94
1	4.9	4.90	4.63
2	4.8	4.80	4.20
3	NaN	NaN	NaN
4	NaN	NaN	NaN

	reviews_per_month	reviews
0	0.52	Barry's place was perfect. It was cute, modern...
1	3.81	I booked this place last minute to attend a fu...
2	2.14	Great spot! Little far out but overall a great...
3	NaN	NaN

[5 rows x 51 columns]

```
In [5]: # missing value proportion
missing_ratio = df.isnull().mean() * 100
print("\nmissing value ratio (%):")
print(missing_ratio)
```

```

missing value ratio (%) :
name                                0.000000
description                         2.465596
property_type                       0.000000
price                               0.000000
neighbourhood_cleansed              0.000000
neighbourhood_group_cleansed        0.000000
latitude                            0.000000
longitude                           0.000000
host_since                          0.000000
host_response_time                  14.035423
host_response_rate                   14.035423
host_acceptance_rate                 13.079766
host_is_superhost                    1.599134
host_listings_count                  0.000000
host_total_listings_count            0.000000
host_verifications                   0.000000
host_has_profile_pic                 0.000000
host_identity_verified               0.000000
calculated_host_listings_count       0.000000
calculated_host_listings_count_entire_homes 0.000000
calculated_host_listings_count_private_rooms 0.000000
calculated_host_listings_count_shared_rooms 0.000000
room_type                           0.000000
accommodates                         0.000000
bathrooms                           0.019113
bathrooms_text                       0.063710
bedrooms                            0.216616
beds                                 0.535168
amenities                            0.000000
has_availability                     0.751784
availability_30                      0.000000
availability_60                      0.000000
availability_90                      0.000000
availability_365                     0.000000
instant_bookable                     0.000000
minimum_nights                       0.000000
maximum_nights                       0.000000
number_of_reviews                    0.000000
number_of_reviews_ltm                 0.000000
number_of_reviews_l30d                0.000000
first_review                         28.504077
last_review                          28.504077
review_scores_rating                  28.504077
review_scores_accuracy                28.504077
review_scores_cleanliness             28.504077
review_scores_checkin                 28.504077
review_scores_communication           28.504077
review_scores_location                28.510449
review_scores_value                   28.504077
reviews_per_month                     28.504077
reviews                              28.504077
dtype: float64

```

```
In [6]: print(df.describe())
```

	price	latitude	longitude	host_response_rate \
count	15696.000000	15696.000000	15696.000000	13493.000000
mean	2.465724	40.726899	-73.943147	91.383013
std	1.709624	0.058079	0.060149	22.320590
min	0.000000	40.500366	-74.251907	0.000000
25%	1.000000	40.685686	-73.983133	97.000000
50%	2.000000	40.725251	-73.952458	100.000000
75%	4.000000	40.762314	-73.921120	100.000000
max	5.000000	40.911390	-73.713650	100.000000

	host_acceptance_rate	host_listings_count	host_total_listings_count
count	13643.000000	15696.000000	15696.000000
mean	78.579198	288.106588	393.435143
std	27.896484	984.327077	1205.427544
min	0.000000	1.000000	1.000000
25%	69.000000	1.000000	2.000000
50%	91.000000	3.000000	5.000000
75%	100.000000	21.000000	31.000000
max	100.000000	4494.000000	9019.000000

	calculated_host_listings_count \
count	15696.000000
mean	74.901631
std	198.042132
min	1.000000
25%	1.000000
50%	3.000000
75%	17.000000
max	876.000000

	calculated_host_listings_count_entire_homes \
count	15696.000000
mean	45.303772
std	166.432525
min	0.000000
25%	0.000000
50%	1.000000
75%	4.000000
max	876.000000

	calculated_host_listings_count_private_rooms ... \
count	15696.000000 ...
mean	27.702281 ...
std	117.504567 ...
min	0.000000 ...
25%	0.000000 ...
50%	1.000000 ...
75%	3.000000 ...
max	719.000000 ...

	number_of_reviews_ltm	number_of_reviews_l30d	review_scores_rating
count	15696.000000	15696.000000	11222.000000
mean	5.693425	0.474134	4.719393
std	23.603555	2.210829	0.462927

min	0.000000	0.000000	1.000000
25%	0.000000	0.000000	4.660000
50%	1.000000	0.000000	4.850000
75%	4.000000	0.000000	5.000000
max	1772.000000	147.000000	5.000000

	review_scores_accuracy	review_scores_cleanliness \
count	11222.000000	11222.000000
mean	4.742812	4.679642
std	0.460347	0.483314
min	1.000000	1.000000
25%	4.690000	4.590000
50%	4.880000	4.820000
75%	5.000000	4.980000
max	5.000000	5.000000

	review_scores_checkin	review_scores_communication \
count	11222.000000	11222.000000
mean	4.82631	4.808233
std	0.37655	0.433165
min	1.00000	1.000000
25%	4.81000	4.800000
50%	4.94000	4.940000
75%	5.00000	5.000000
max	5.00000	5.000000

	review_scores_location	review_scores_value	reviews_per_month
count	11221.000000	11222.000000	11222.000000
mean	4.721844	4.609505	1.245801
std	0.400359	0.512808	2.269312
min	1.000000	1.000000	0.010000
25%	4.630000	4.510000	0.210000
50%	4.820000	4.750000	0.610000
75%	5.000000	4.890000	1.650000
max	5.000000	5.000000	110.100000

[8 rows x 32 columns]

```
In [7]: categorical_columns = df.select_dtypes(include=['object']).columns
for col in categorical_columns:
    print(f"column_name: {col}, unique value: {df[col].nunique()}")
```

column_name: name, unique value: 15189
column_name: description, unique value: 12687
column_name: property_type, unique value: 59
column_name: neighbourhood_cleansed, unique value: 217
column_name: neighbourhood_group_cleansed, unique value: 5
column_name: host_since, unique value: 4037
column_name: host_response_time, unique value: 4
column_name: host_is_superhost, unique value: 2
column_name: host_verifications, unique value: 6
column_name: room_type, unique value: 4
column_name: bathrooms_text, unique value: 30
column_name: amenities, unique value: 13314
column_name: has_availability, unique value: 1
column_name: first_review, unique value: 3261
column_name: last_review, unique value: 1390
column_name: reviews, unique value: 11215

```
In [8]: columns_to_drop = ['name', 'description', 'reviews', 'amenities', 'neighbourh  
df = df.drop(columns=[col for col in columns_to_drop if col in df.columns])  
  
print(df.head())
```

	property_type	price	neighbourhood_group_cleansed	latitude
0	Entire rental unit	4	Brooklyn	40.68456
1	Private room in rental unit	3	Brooklyn	40.63899
2	Entire rental unit	3	Brooklyn	40.61881
3	Private room in rental unit	0	Brooklyn	40.67397
4	Room in hotel	2	Manhattan	40.74718

	longitude	host_since	host_response_time	host_response_rate
0	-73.939870	2015-05-23 00:00:00	within a day	100.0
1	-73.965739	2023-09-14 00:00:00	within an hour	100.0
2	-74.032380	2022-07-31 00:00:00	within an hour	100.0
3	-73.953990	2012-08-11 00:00:00	within an hour	99.0
4	-73.985390	2014-12-23 00:00:00	within an hour	93.0

	host_acceptance_rate	host_is_superhost	...	first_review
0	100.0	True	...	2019-04-28 00:00:00
1	98.0	True	...	2024-01-13 00:00:00
2	100.0	False	...	2024-06-27 00:00:00
3	23.0	False	...	NaN
4	95.0	False	...	NaN

	last_review	review_scores_rating	review_scores_accuracy
0	2024-08-10 00:00:00	5.00	5.00
1	2024-09-02 00:00:00	4.83	4.87
2	2024-08-17 00:00:00	4.60	4.80
3	NaN	NaN	NaN
4	NaN	NaN	NaN

	review_scores_cleanliness	review_scores_checkin
0	4.97	5.0
1	4.93	4.8
2	4.20	4.8
3	NaN	NaN
4	NaN	NaN

	review_scores_communication	review_scores_location	review_scores_value
0	5.0	4.71	4.94
1	4.9	4.90	4.63
2	4.8	4.80	4.20
3	NaN	NaN	NaN
4	NaN	NaN	NaN

	reviews_per_month
0	0.52
1	3.81
2	2.14
3	NaN
4	NaN

[5 rows x 44 columns]

```
In [9]: df['host_since'] = pd.to_datetime(df['host_since'], errors='coerce')
df['host_since_year_diff'] = 2024 - df['host_since'].dt.year
print(df[['host_since', 'host_since_year_diff']].head())
```

	host_since	host_since_year_diff
0	2015-05-23	9
1	2023-09-14	1
2	2022-07-31	2
3	2012-08-11	12
4	2014-12-23	10

```
In [10]: df['last_review'] = pd.to_datetime(df['last_review'], errors='coerce')
df['first_review'] = pd.to_datetime(df['first_review'], errors='coerce')

df['days_since_last_review'] = (pd.to_datetime('today') - df['last_review']).dt.days
df['days_since_first_review'] = (pd.to_datetime('today') - df['first_review']).dt.days

def categorize_review_period(days):
    if days <= 365:
        return 'within 1 year'
    elif days <= 1095:
        return '1-3 years'
    else:
        return 'more than 3 years'

df['review_last_category'] = df['days_since_last_review'].apply(categorize_review_period)
df['review_first_category'] = df['days_since_first_review'].apply(categorize_review_period)

print(df[['last_review', 'days_since_last_review', 'review_last_category']].head())
print(df[['first_review', 'days_since_first_review', 'review_first_category']].head())
```

	last_review	days_since_last_review	review_last_category
0	2024-08-10	105.0	within 1 year
1	2024-09-02	82.0	within 1 year
2	2024-08-17	98.0	within 1 year
3	NaT	NaN	more than 3 years
4	NaT	NaN	more than 3 years

	first_review	days_since_first_review	review_first_category
0	2019-04-28	2036.0	more than 3 years
1	2024-01-13	315.0	within 1 year
2	2024-06-27	149.0	within 1 year
3	NaT	NaN	more than 3 years
4	NaT	NaN	more than 3 years

```
In [11]: df.head(5)
```

```
Out[11]:
```

	property_type	price	neighbourhood_group_cleansed	latitude	longitude
0	Entire rental unit	4	Brooklyn	40.684560	-73.939870
1	Private room in rental unit	3	Brooklyn	40.638991	-73.965739
2	Entire rental unit	3	Brooklyn	40.618810	-74.032380
3	Private room in rental unit	0	Brooklyn	40.673970	-73.953990
4	Room in hotel	2	Manhattan	40.747180	-73.985390

5 rows × 49 columns

```
In [12]: room_type_counts = df['room_type'].value_counts()

print(room_type_counts)
```

```
room_type
Entire home/apt      8592
Private room         6737
Hotel room           200
Shared room          167
Name: count, dtype: int64
```

```
In [13]: neighbourhood_to_score = {
    'Manhattan': 5,
    'Brooklyn': 4,
    'Queens': 3,
    'Bronx': 2,
    'Staten Island': 1
}

df['neighbourhood_encoded'] = df['neighbourhood_group_cleansed'].map(neighbourhood_to_score)
df.head(5)
```

```
Out[13]:
```

	property_type	price	neighbourhood_group_cleansed	latitude	longitude
0	Entire rental unit	4	Brooklyn	40.684560	-73.939870
1	Private room in rental unit	3	Brooklyn	40.638991	-73.965739
2	Entire rental unit	3	Brooklyn	40.618810	-74.032380
3	Private room in rental unit	0	Brooklyn	40.673970	-73.953990
4	Room in hotel	2	Manhattan	40.747180	-73.985390

5 rows × 50 columns

```
In [14]: response_counts = df['host_response_time'].value_counts()
last_counts = df['review_last_category'].value_counts()
first_counts = df['review_first_category'].value_counts()

print(response_counts)
print(last_counts)
print(first_counts)
```

```
host_response_time
within an hour      8367
within a few hours  2908
within a day        1472
a few days or more   746
Name: count, dtype: int64
review_last_category
within 1 year      6972
more than 3 years  5204
1-3 years          3520
Name: count, dtype: int64
review_first_category
more than 3 years   9528
1-3 years          4397
within 1 year       1771
Name: count, dtype: int64
```

```
In [15]: response_to_score = {
    'within an hour': 4,
    'within a few hours': 3,
    'within a day': 2,
    'a few days or more': 1,
}

df['response_encoded'] = df['host_response_time'].map(response_to_score)

last_to_score = {
    'within 1 year': 3,
    '1-3 years': 2,
    'more than 3 years': 1,
}
df['last_encoded'] = df['review_last_category'].map(last_to_score)

df_1 = pd.get_dummies(df, columns=['review_first_category', 'room_type'],
                      prefix=['first_review', 'room_type'])

df_1.head(5)
```

Out[15]:

	property_type	price	neighbourhood_group_cleansed	latitude	longitude
0	Entire rental unit	4	Brooklyn	40.684560	-73.939870
1	Private room in rental unit	3	Brooklyn	40.638991	-73.965739
2	Entire rental unit	3	Brooklyn	40.618810	-74.032380
3	Private room in rental unit	0	Brooklyn	40.673970	-73.953990
4	Room in hotel	2	Manhattan	40.747180	-73.985390

5 rows × 57 columns

In [16]: `df = df_1.drop(columns=['property_type', 'neighbourhood_group_cleansed', 'host'`

```
all_features = df.columns
print(all_features)
```

```
Index(['price', 'latitude', 'longitude', 'host_response_rate',
      'host_acceptance_rate', 'host_is_superhost', 'host_listings_count',
      'host_total_listings_count', 'host_has_profile_pic',
      'host_identity_verified', 'calculated_host_listings_count',
      'calculated_host_listings_count_entire_homes',
      'calculated_host_listings_count_private_rooms',
      'calculated_host_listings_count_shared_rooms', 'accommodates',
      'bathrooms', 'bedrooms', 'beds', 'has_availability', 'availability_30',
      'availability_60', 'availability_90', 'availability_365',
      'instant_bookable', 'minimum_nights', 'maximum_nights',
      'number_of_reviews', 'number_of_reviews_ltm', 'number_of_reviews_l30d',
      'review_scores_rating', 'review_scores_accuracy',
      'review_scores_cleanliness', 'review_scores_checkin',
      'review_scores_communication', 'review_scores_location',
      'review_scores_value', 'reviews_per_month', 'host_since_year_diff',
      'days_since_last_review', 'days_since_first_review',
      'neighbourhood_encoded', 'response_encoded', 'last_encoded',
      'first_review_1-3 years', 'first_review_more than 3 years',
      'first_review_within 1 year', 'room_type_Entire home/apt',
      'room_type_Hotel room', 'room_type_Private room',
      'room_type_Shared room'],
      dtype='object')
```

In [17]: `# check missing value
print("\nmissing value stats:")
print(df.isnull().sum())`

```
# visualization
```

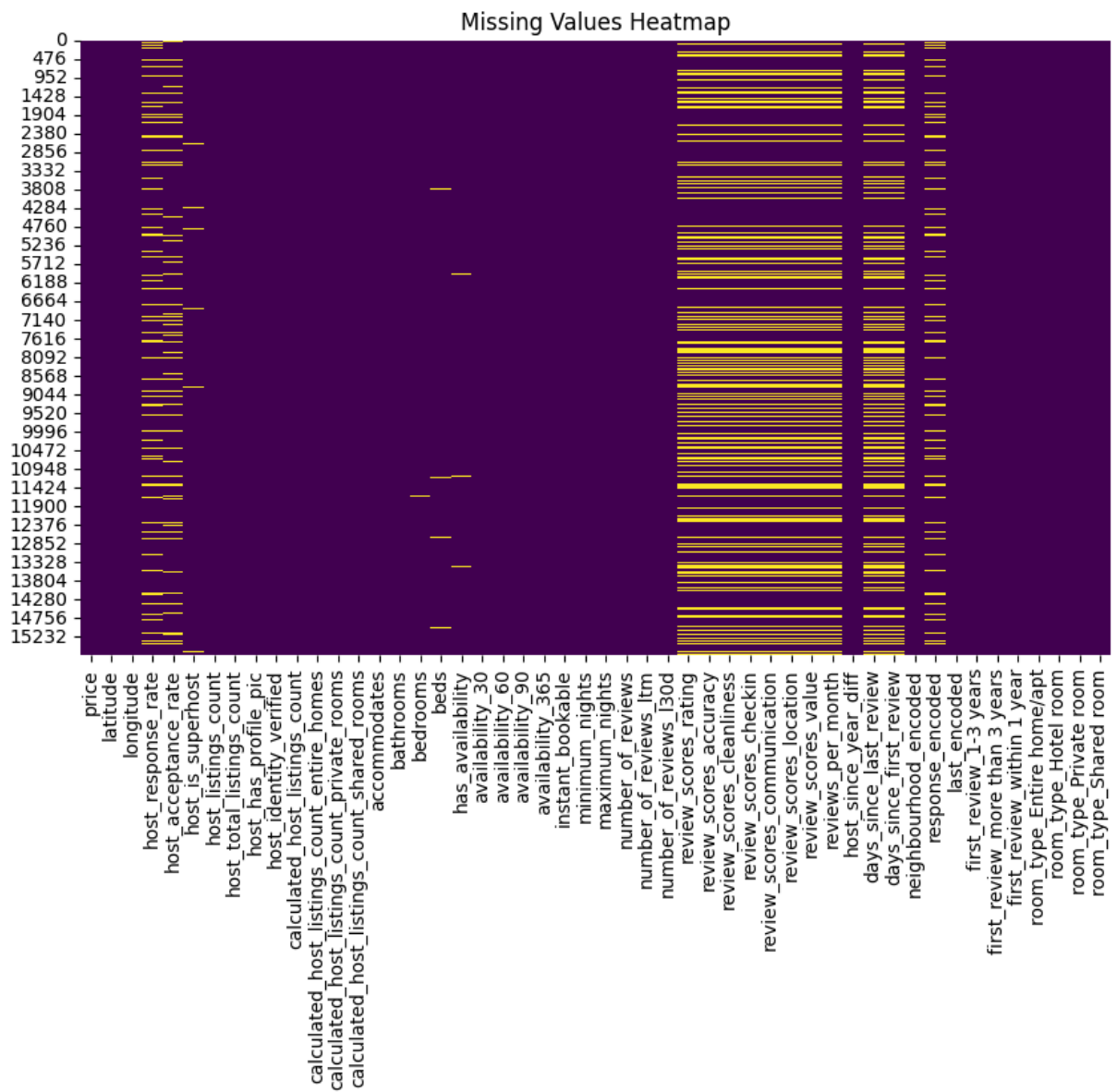
```
plt.figure(figsize=(10, 6))
sns.heatmap(df.isnull(), cbar=False, cmap="viridis")
```

```
plt.title("Missing Values Heatmap")
plt.show()
```

missing value stats:

price	0
latitude	0
longitude	0
host_response_rate	2203
host_acceptance_rate	2053
host_is_superhost	251
host_listings_count	0
host_total_listings_count	0
host_has_profile_pic	0
host_identity_verified	0
calculated_host_listings_count	0
calculated_host_listings_count_entire_homes	0
calculated_host_listings_count_private_rooms	0
calculated_host_listings_count_shared_rooms	0
accommodates	0
bathrooms	3
bedrooms	34
beds	84
has_availability	118
availability_30	0
availability_60	0
availability_90	0
availability_365	0
instant_bookable	0
minimum_nights	0
maximum_nights	0
number_of_reviews	0
number_of_reviews_ltm	0
number_of_reviews_l30d	0
review_scores_rating	4474
review_scores_accuracy	4474
review_scores_cleanliness	4474
review_scores_checkin	4474
review_scores_communication	4474
review_scores_location	4475
review_scores_value	4474
reviews_per_month	4474
host_since_year_diff	0
days_since_last_review	4474
days_since_first_review	4474
neighbourhood_encoded	0
response_encoded	2203
last_encoded	0
first_review_1-3 years	0
first_review_more than 3 years	0
first_review_within 1 year	0
room_type_Entire home/apt	0
room_type_Hotel room	0
room_type_Private room	0
room_type_Shared room	0

dtype: int64



```
In [18]: # missing value proportion
missing_ratio = df.isnull().mean() * 100
print("\nmissing value ratio (%) :")
print(missing_ratio)
```

```

missing value ratio (%) :
price                0.000000
latitude             0.000000
longitude            0.000000
host_response_rate   14.035423
host_acceptance_rate 13.079766
host_is_superhost     1.599134
host_listings_count   0.000000
host_total_listings_count 0.000000
host_has_profile_pic  0.000000
host_identity_verified 0.000000
calculated_host_listings_count 0.000000
calculated_host_listings_count_entire_homes 0.000000
calculated_host_listings_count_private_rooms 0.000000
calculated_host_listings_count_shared_rooms 0.000000
accommodates         0.000000
bathrooms            0.019113
bedrooms             0.216616
beds                 0.535168
has_availability      0.751784
availability_30       0.000000
availability_60       0.000000
availability_90       0.000000
availability_365      0.000000
instant_bookable      0.000000
minimum_nights        0.000000
maximum_nights        0.000000
number_of_reviews     0.000000
number_of_reviews_ltm 0.000000
number_of_reviews_l30d 0.000000
review_scores_rating   28.504077
review_scores_accuracy 28.504077
review_scores_cleanliness 28.504077
review_scores_checkin  28.504077
review_scores_communication 28.504077
review_scores_location 28.510449
review_scores_value    28.504077
reviews_per_month     28.504077
host_since_year_diff   0.000000
days_since_last_review 28.504077
days_since_first_review 28.504077
neighbourhood_encoded  0.000000
response_encoded       14.035423
last_encoded           0.000000
first_review_1-3 years 0.000000
first_review_more than 3 years 0.000000
first_review_within 1 year 0.000000
room_type_Entire home/apt 0.000000
room_type_Hotel room    0.000000
room_type_Private room  0.000000
room_type_Shared room   0.000000
dtype: float64

```

```

In [19]: # 找到缺失值比例大于 90% 的列并打印这些列名
columns_with_high_missing = missing_ratio[missing_ratio > 0.9].index

```

```

if len(columns_with_high_missing) > 0:
    print("Columns with missing value ratio > 90%:")
    print(columns_with_high_missing)
else:
    print("No columns have missing value ratio greater than 90%.")

columns_with_high_missing = missing_ratio[missing_ratio > 0.7].index

if len(columns_with_high_missing) > 0:
    print("Columns with missing value ratio > 70%:")
    print(columns_with_high_missing)
else:
    print("No columns have missing value ratio greater than 70%.")

```

Columns with missing value ratio > 90%:

```

Index(['host_response_rate', 'host_acceptance_rate', 'host_is_superhost',
      'review_scores_rating', 'review_scores_accuracy',
      'review_scores_cleanliness', 'review_scores_checkin',
      'review_scores_communication', 'review_scores_location',
      'review_scores_value', 'reviews_per_month', 'days_since_last_review',
      'days_since_first_review', 'response_encoded'],
      dtype='object')

```

Columns with missing value ratio > 70%:

```

Index(['host_response_rate', 'host_acceptance_rate', 'host_is_superhost',
      'has_availability', 'review_scores_rating', 'review_scores_accuracy',
      'review_scores_cleanliness', 'review_scores_checkin',
      'review_scores_communication', 'review_scores_location',
      'review_scores_value', 'reviews_per_month', 'days_since_last_review',
      'days_since_first_review', 'response_encoded'],
      dtype='object')

```

In [20]: `df = df.drop(columns=['has_availability', 'reviews_per_month', 'days_since_la`

```

In [22]: # missing value proportion
missing_ratio = df.isnull().mean() * 100
print("\nmissing value ratio (%) :")
print(missing_ratio)

```

missing value ratio (%) :	
price	0.000000
latitude	0.000000
longitude	0.000000
host_response_rate	0.000000
host_acceptance_rate	0.000000
host_is_superhost	1.599134
host_listings_count	0.000000
host_total_listings_count	0.000000
host_has_profile_pic	0.000000
host_identity_verified	0.000000
calculated_host_listings_count	0.000000
calculated_host_listings_count_entire_homes	0.000000
calculated_host_listings_count_private_rooms	0.000000
calculated_host_listings_count_shared_rooms	0.000000
accommodates	0.000000
bathrooms	0.019113
bedrooms	0.216616
beds	0.535168
availability_30	0.000000
availability_60	0.000000
availability_90	0.000000
availability_365	0.000000
instant_bookable	0.000000
minimum_nights	0.000000
maximum_nights	0.000000
number_of_reviews	0.000000
number_of_reviews_ltm	0.000000
number_of_reviews_l30d	0.000000
review_scores_rating	0.000000
review_scores_accuracy	0.000000
review_scores_cleanliness	0.000000
review_scores_checkin	0.000000
review_scores_communication	0.000000
review_scores_location	0.000000
review_scores_value	0.000000
host_since_year_diff	0.000000
neighbourhood_encoded	0.000000
last_encoded	0.000000
first_review_1-3 years	0.000000
first_review_more than 3 years	0.000000
first_review_within 1 year	0.000000
room_type_Entire home/apt	0.000000
room_type_Hotel room	0.000000
room_type_Private room	0.000000
room_type_Shared room	0.000000
dtype: float64	

```
In [23]: missing_count = df.isnull().sum()
print("Missing values count in each column:")
missing_count
```

Missing values count in each column:

```

Out[23]: price                                0
latitude                                     0
longitude                                    0
host_response_rate                          0
host_acceptance_rate                       0
host_is_superhost                           251
host_listings_count                         0
host_total_listings_count                  0
host_has_profile_pic                        0
host_identity_verified                     0
calculated_host_listings_count              0
calculated_host_listings_count_entire_homes 0
calculated_host_listings_count_private_rooms 0
calculated_host_listings_count_shared_rooms 0
accommodates                               0
bathrooms                                  3
bedrooms                                   34
beds                                       84
availability_30                            0
availability_60                            0
availability_90                            0
availability_365                           0
instant_bookable                           0
minimum_nights                             0
maximum_nights                             0
number_of_reviews                          0
number_of_reviews_ltm                      0
number_of_reviews_l30d                    0
review_scores_rating                       0
review_scores_accuracy                     0
review_scores_cleanliness                  0
review_scores_checkin                     0
review_scores_communication                0
review_scores_location                     0
review_scores_value                        0
host_since_year_diff                       0
neighbourhood_encoded                      0
last_encoded                              0
first_review_1-3 years                     0
first_review_more than 3 years              0
first_review_within 1 year                 0
room_type_Entire home/apt                  0
room_type_Hotel room                       0
room_type_Private room                     0
room_type_Shared room                      0
dtype: int64

```

```

In [25]: import matplotlib.pyplot as plt

variables = ['accommodates', 'bathrooms', 'bedrooms', 'beds']

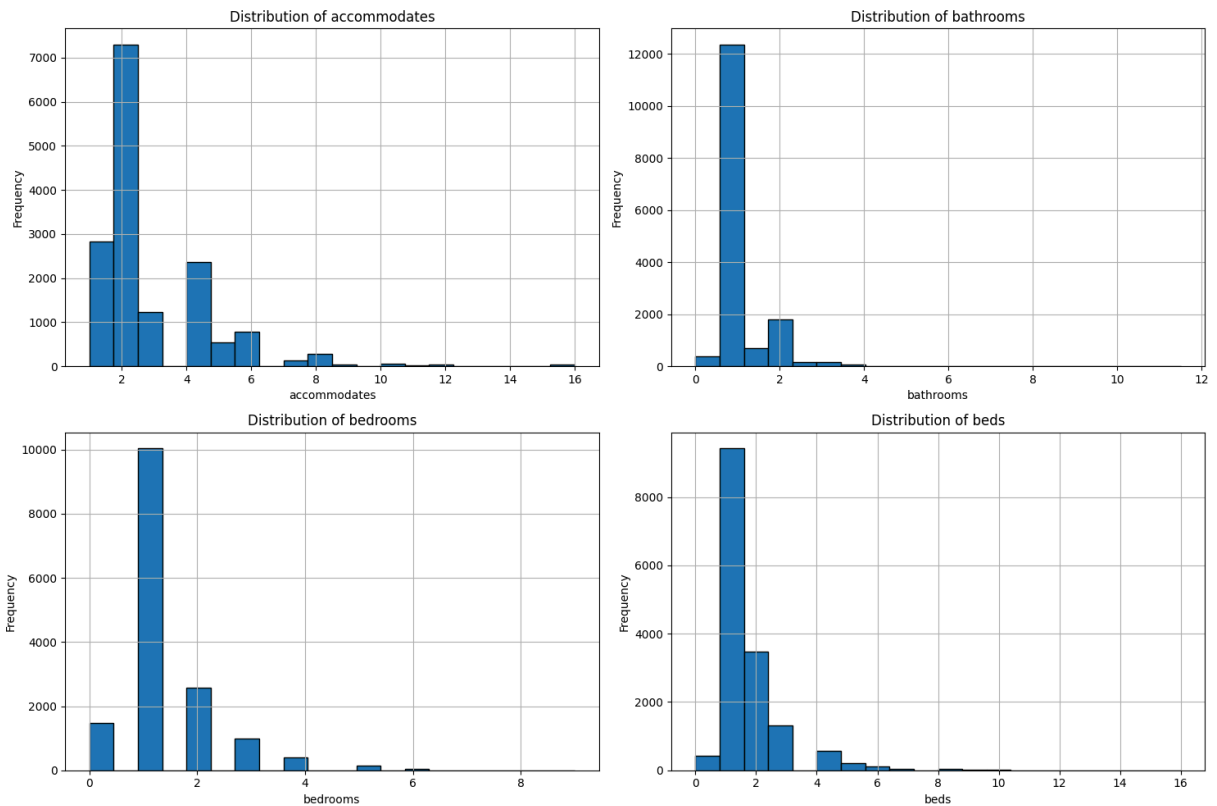
plt.figure(figsize=(15, 10))

for i, var in enumerate(variables, 1):
    plt.subplot(2, 2, i)
    df[var].hist(bins=20, edgecolor='black')

```

```
plt.xlabel(var)
plt.ylabel('Frequency')
plt.title(f'Distribution of {var}')

plt.tight_layout()
plt.show()
```



```
In [27]: missing_variables = df.columns[df.isnull().any()]
print("Variables with missing values:")
print(missing_variables)
```

Variables with missing values:
Index([], dtype='object')

feature importance

```
In [28]: from sklearn.ensemble import RandomForestRegressor
import pandas as pd
import matplotlib.pyplot as plt

X = df.drop(columns=['price'])
y = df['price']

model = RandomForestRegressor(n_estimators=100, random_state=100)
model.fit(X, y)

feature_importances = model.feature_importances_
```

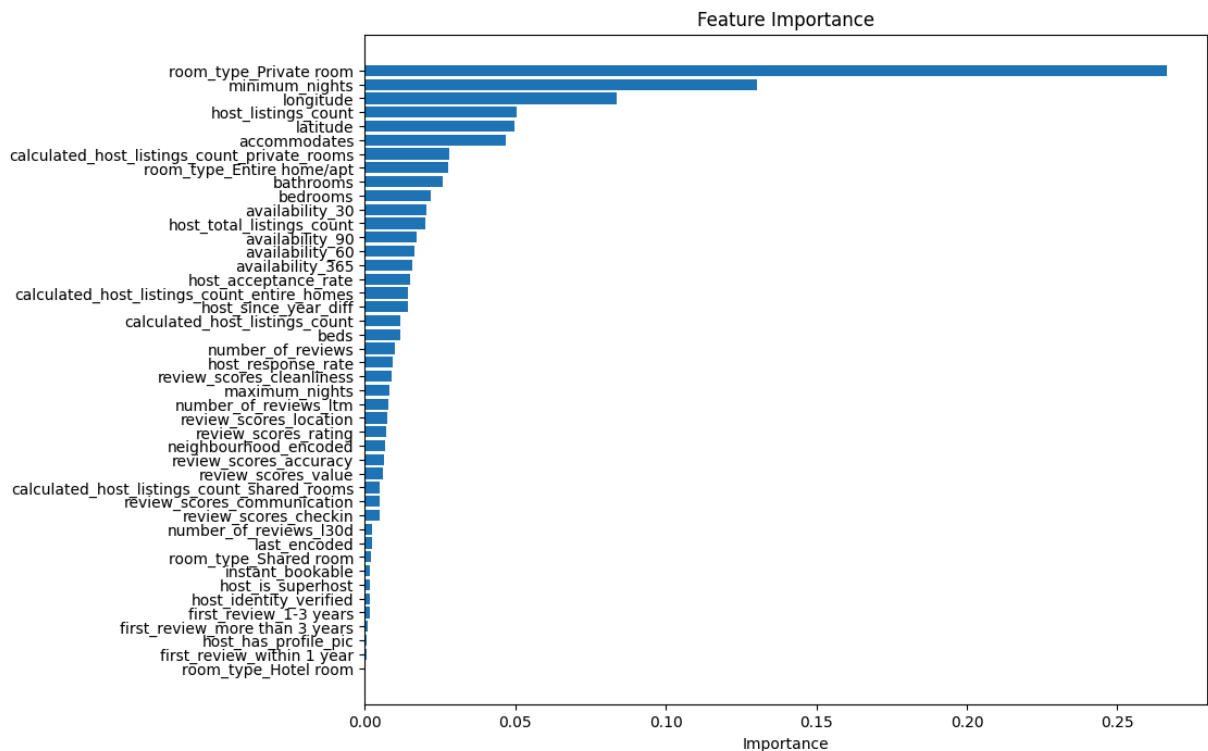
```
importance_df = pd.DataFrame({
    'Feature': X.columns,
    'Importance': feature_importances
}).sort_values(by='Importance', ascending=False)

print("Feature Importance:")
print(importance_df)

plt.figure(figsize=(10, 8))
plt.barh(importance_df['Feature'], importance_df['Importance'])
plt.xlabel('Importance')
plt.title('Feature Importance')
plt.gca().invert_yaxis()
plt.show()
```

Feature Importance:

	Feature	Importance
42	room_type_Private room	0.266504
22	minimum_nights	0.130151
1	longitude	0.083623
5	host_listings_count	0.050381
0	latitude	0.049821
13	accommodates	0.046973
11	calculated_host_listings_count_private_rooms	0.028074
40	room_type_Entire home/apt	0.027733
14	bathrooms	0.025866
15	bedrooms	0.021874
17	availability_30	0.020384
6	host_total_listings_count	0.020270
19	availability_90	0.017434
18	availability_60	0.016413
20	availability_365	0.015812
3	host_acceptance_rate	0.015218
10	calculated_host_listings_count_entire_homes	0.014353
34	host_since_year_diff	0.014282
9	calculated_host_listings_count	0.011861
16	beds	0.011782
24	number_of_reviews	0.010107
2	host_response_rate	0.009216
29	review_scores_cleanliness	0.009097
23	maximum_nights	0.008237
25	number_of_reviews_ltm	0.007911
32	review_scores_location	0.007551
27	review_scores_rating	0.007337
35	neighbourhood_encoded	0.006751
28	review_scores_accuracy	0.006318
33	review_scores_value	0.006192
12	calculated_host_listings_count_shared_rooms	0.005187
31	review_scores_communication	0.005028
30	review_scores_checkin	0.004958
26	number_of_reviews_l30d	0.002595
36	last_encoded	0.002453
43	room_type_Shared room	0.002180
21	instant_bookable	0.001926
4	host_is_superhost	0.001793
8	host_identity_verified	0.001752
37	first_review_1-3 years	0.001656
38	first_review_more than 3 years	0.001071
7	host_has_profile_pic	0.000861
39	first_review_within 1 year	0.000687
41	room_type_Hotel room	0.000325



```
In [ ]: from catboost import CatBoostClassifier

X = df.drop(columns=['price'])
y = df['price']

X_train, X_test, y_train, y_test = df(X, y, test_size=0.2, random_state=1005)

model_cb = CatBoostClassifier(iterations=1000, learning_rate=0.1, depth=10,
model_cb.fit(X_train, y_train)

y_pred = model_cb.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
print("\nClassification Report:\n", classification_report(y_test, y_pred))
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print("\nRoot Mean Squared Error (RMSE):", rmse)
```

```
In [ ]: import xgboost as xgb
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.metrics import mean_squared_error, r2_score

X = df.drop(columns=['price'])
y = df['price']
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=1005)

param_grid = {
    'n_estimators': [100, 150, 200],
    'learning_rate': [0.01, 0.05, 0.1],
    'max_depth': [4, 5, 6]
}
```

```

xgb_model = xgb.XGBRegressor(objective='reg:squarederror', random_state=1005)

grid_search = GridSearchCV(estimator=xgb_model, param_grid=param_grid, cv=5,

grid_search.fit(X_train, y_train)

best_model = grid_search.best_estimator_
print("Best Parameters:", grid_search.best_params_)

y_pred = best_model.predict(X_val)

mse = mean_squared_error(y_val, y_pred)
r2 = r2_score(y_val, y_pred)
print(f"Mean Squared Error: {mse}")
print(f"R^2 Score: {r2}")

file_path_1 = '/home/users/ys468/ml_kaggle/cs-671-fall-2024-final-project/te
df_t = pd.read_csv(file_path_1)

categorical_columns = df_t.select_dtypes(include=['object']).columns
for col in categorical_columns:
    print(f"column_name: {col}, unique value: {df_t[col].nunique()}")

columns_to_drop = ['name', 'description', 'reviews', 'amenities', 'neighbourh
df_t = df_t.drop(columns=[col for col in columns_to_drop if col in df_t.colu

print(df_t.head())

df_t['host_since'] = pd.to_datetime(df_t['host_since'], errors='coerce')
df_t['host_since_year_diff'] = 2024 - df_t['host_since'].dt.year
print(df_t[['host_since', 'host_since_year_diff']].head())

df_t['last_review'] = pd.to_datetime(df_t['last_review'], errors='coerce')
df_t['first_review'] = pd.to_datetime(df_t['first_review'], errors='coerce')

df_t['days_since_last_review'] = (pd.to_datetime('today') - df_t['last_review']).dt.days
df_t['days_since_first_review'] = (pd.to_datetime('today') - df_t['first_review']).dt.days

def categorize_review_period(days):
    if days <= 365:
        return 'within 1 year'
    elif days <= 1095:
        return '1-3 years'
    else:
        return 'more than 3 years'

```

```

df_t['review_last_category'] = df_t['days_since_last_review'].apply(categori
df_t['review_first_category'] = df_t['days_since_first_review'].apply(categori

print(df_t[['last_review', 'days_since_last_review', 'review_last_category']]
print(df_t[['first_review', 'days_since_first_review', 'review_first_category']]

room_type_counts = df_t['room_type'].value_counts()

print(room_type_counts)

neighbourhood_to_score = {
    'Manhattan': 5,
    'Brooklyn': 4,
    'Queens': 3,
    'Bronx': 2,
    'Staten Island': 1
}

df_t['neighbourhood_encoded'] = df_t['neighbourhood_group_cleansed'].map(nei

response_counts = df_t['host_response_time'].value_counts()
last_counts = df_t['review_last_category'].value_counts()
first_counts = df_t['review_first_category'].value_counts()

print(response_counts)
print(last_counts)
print(first_counts)

response_to_score = {
    'within an hour': 4,
    'within a few hours': 3,
    'within a day': 2,
    'a few days or more': 1,
}

df_t['response_encoded'] = df_t['host_response_time'].map(response_to_score)

last_to_score = {
    'within 1 year': 3,
    '1-3 years': 2,
    'more than 3 years': 1,
}
df_t['last_encoded'] = df_t['review_last_category'].map(last_to_score)

df_t_1 = pd.get_dummies(df_t, columns=['review_first_category', 'room_type'],
                        prefix=['first_review', 'room_type'])

df_t = df_t_1.drop(columns=['property_type', 'neighbourhood_group_cleansed'],

all_features = df_t.columns
print(all_features)

```

```

df_t = df_t.drop(columns=[ 'id','has_availability','reviews_per_month','days

columns_to_fill = ['host_response_rate', 'host_acceptance_rate', 'review_sco
                    'review_scores_accuracy', 'review_scores_cleanliness',
                    'review_scores_checkin', 'review_scores_communication',
                    'review_scores_location', 'review_scores_value']

for column in columns_to_fill:
    mode_value = df_t[column].mode()[0]
    df_t[column].fillna(mode_value, inplace=True)

mode_value = df_t['host_is_superhost'].mode()[0]
df_t['host_is_superhost'].fillna(mode_value, inplace=True)

variables = ['accommodates', 'bathrooms', 'bedrooms', 'beds']

plt.figure(figsize=(15, 10))

for i, var in enumerate(variables, 1):
    plt.subplot(2, 2, i)
    df_t[var].hist(bins=20, edgecolor='black')
    plt.xlabel(var)
    plt.ylabel('Frequency')
    plt.title(f'Distribution of {var}')

plt.tight_layout()
plt.show()

columns_to_fill_mode = ['accommodates']
columns_to_fill_median = ['bathrooms', 'bedrooms', 'beds']

for column in columns_to_fill_mode:
    mode_value = df_t[column].mode()[0]
    df_t[column].fillna(mode_value, inplace=True)

for column in columns_to_fill_median:
    median_value = df_t[column].median()
    df_t[column].fillna(median_value, inplace=True)

X_test = df_t
y_pred = best_model.predict(X_test)
df_t['price'] = y_pred

df_t.to_csv('test_with_predictions.csv', index=False)

print("Predictions saved to 'test_with_predictions.csv'")

```