AIRBNB RENT PREDICTION MODEL BASED ON LARGE DATA

A SENIOR PROJECT

SUBMITTED TO THE COLLEGE OF ENGINEERING

OF TENNESSEE STATE UNIVERSITY

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

BACHELOR OF SCIENCE IN COMPUTER SCIENCE

Eulice Winston IV

Nicholas Stepka

DECEMBER 2023

To the Department of Computer Science,

We are submitting a senior project by "Eulice Winston IV" and "Nicholas Stepka" entitled "Airbnb Rent Prediction Model Based in Large Data". We recommend that it be accepted in partial fulfillment of the requirements for the degree, Bachelor of Science in Computer Science.

 

 

Technical Advisor, Dr. Manar Samad

 

Course Instructor, Dr. Ali Sekmen

 

Department Chair, Dr. Tamara Rogers

 

Accepted for the College of Engineering:

Dr. Lin Li

 

Dean of the College of Engineering

# ABSTRACT

This study explores the application of machine learning algorithms to establish fair rental prices in the dynamic real estate market, focusing on Airbnb listings in Nashville. Given the challenge of determining reasonable rental prices for both owners and renters, we employ a comprehensive approach combining data preprocessing, advanced regression modeling, and causal inference. Our methodology includes analyzing a large dataset using various machine learning models, with the HistGradientBoosting Regressor identified as the most effective for predicting rental prices. We emphasize feature importance and reduction to enhance computational efficiency while maintaining statistical accuracy.

The research further delves into causal analysis, utilizing state-of-the-art algorithms and domain expertise to discern the impact of specific property attributes on rental pricing, particularly examining the role of 'accommodates' in influencing price. The study employs rigorous empirical testing, including refutation methods, to validate the causal relationships and ensure the robustness of our findings.

The culmination of this project is a user-friendly web application developed using Python Flask, enabling users to input property features and receive predicted rental prices along with relevant metrics. This tool serves as a practical solution for property owners, renters, and investors in the Airbnb market, aiding in informed decision-making. Our study demonstrates the effectiveness of integrating machine learning, causal inference, and domain knowledge to address complex challenges in the real estate sector.

In terms of results, our HistGradientBoosting regression model achieved a commendable R2 score of 0.70, significantly enhancing predictive accuracy by reducing the feature set from over 200 to just 20. Additionally, for our causal question, we obtained a P-score of 0.98, affirming that the 'accommodates' feature significantly influences the pricing of Airbnb listings.

# ACKNOWLEDGMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LISTINGS

# CHAPTER 1

# INTRODUCTION

This study aims to investigate the feasibility of using machine learning algorithms to establish fair rental prices in the constantly evolving real estate market. The challenge lies in the difficulty for both owners and renters to agree on reasonable rental prices, leading to missed opportunities and incorrect pricing decisions.

To address this issue, we plan to employ machine learning algorithms to analyze a large rental data set and develop a fair pricing model. The model will consider several key variables, including demographic information and property characteristics such as location, size, and amenities. By considering these variables, the model will be able to provide an accurate and scientifically validated market price for rental properties, allowing for confident rental decisions for both owners and renters.

In addition, the model will provide a quantitative importance score for each variable, allowing for further analysis and understanding of the factors that influence rental pricing. This tool will also benefit real estate investors, as they can use the model to determine the potential profitability of a property.

By utilizing state-of-the-art natural language processing techniques and machine learning algorithms, this project aims to provide a scientifically validated solution to the issue of determining fair rental prices in the constantly changing real estate market.

## 1.1     Need Analysis

The rental market is complex and constantly evolving, making it difficult for both renters and owners to agree on fair rental prices. There is a need for a scientifically validated pricing system that considers multiple variables such as demographic data, location, size, and amenities to provide

an accurate and competitive market price for rental properties. Additionally, there is a need for a tool that can provide a quantitative importance score for each variable to help renters and owners make more informed rental decisions. The target audience for this tool includes renters, owners, and investors in the real estate market. Renters and owners will benefit from the tool as it will help them make informed rental decisions by providing them with accurate market prices and the importance score of each variable. Investors will be able to use the tool to determine the potential profitability of a property, helping them make better investment decisions.

## 1.2    Problem Statement

The rental market is constantly changing due to inflation. It is challenging for owners and renters to agree on fair rental prices. In the short-term, indecision leads to missed opportunities and overpaying for rent. While for owners pricing too high can lead to additional months where the property is not rented, a lower price may miss the opportunity for increased profits. This project aims to develop a fair pricing model by mining a large rental data set using machine learning (ML) algorithms. The model will consider many variables, including demographic data, to determine a competitive market price for rental properties. The model will analyze variables such as location, size, and amenities to provide accurate pricing, giving both owners and renters confidence in their rental decisions. This will additionally provide the quantitative importance score for individual variables in making the rental decision. In addition, investors in the real estate market will be able to use the tool to determine if a property on the market has the potential to have a profitable return on investment. The project will entail using state-of-the-art natural language processing (NLP) and ML techniques, providing a scientifically validated pricing system for the housing market.

## 1.3    Literature Review

The sharing economy has disrupted the traditional hospitality industry, and the rise of platforms such as Airbnb has created opportunities for property owners to rent their spaces and for travelers

to experience unique and affordable accommodations. However, determining fair and competitive rental prices is a complex and challenging issue due to the lack of industry standards and the growing number of listings. To address this issue, machine learning (ML) and natural language processing (NLP) have emerged as powerful tools in the sharing economy to analyze data and improve pricing models.

### 1.3.1    Tabular Data and Airbnb Pricing

Tabular data, including tables, has been widely used in analyzing Airbnb pricing models. By leveraging tabular data, analysts can consider various factors such as location, amenities, and seasonality to develop pricing models for Airbnb rentals. Located on Kaggle [1], we will use a big data set with over 134,000 listings with over 100 variables. Through tabular data analysis, rental prices can be predicted accurately, helping property owners to set optimal prices and maximize profits.

### 1.3.2    Regression Model and Airbnb Pricing

Regression analysis is a statistical technique commonly used to analyze the relationship between variables in order to cause predictions [2]. By utilizing machine learning models such as linear regression, gradient boosting regression, and random Forest regression, analysts can predict rental prices accurately and identify the most significant factors that influence rental prices [5]. Furthermore, regularization and feature selection techniques can help avoid overfitting and select the most relevant features. This section will review the existing literature on the application of machine learning models for regression analysis in Airbnb pricing.

### 1.3.3    Different Machine Learning Models for Airbnb Pricing

According to an article by Garcia and peers, "machine learning algorithms (ML) are increasingly being used for the mass appraisal of real estate and in automated valuation models successfully." [8] Such models as random forest, gradient boosting, and neural networks, have been employed

to predict the real estate market. In particular, a convolutional neural network has been utilized to successfully forecast rental prices based on images and text data. In an interview given to Venturebeat.com, Mike Curtis, the VP of Engineering at Airbnb, said they have implemented a machine-learned ranking model to offer predictive rankings of properties [9]. The available literature suggests that ML and natural language processing (NLP) techniques can be utilized to develop accurate and personalized pricing models for short-term rentals. Furthermore, regression analysis and tabular data can also be utilized to explore the factors that influence rental prices. Further research is required to explore the potential of different ML models and to develop more precise and personalized pricing models for the real-estate market.

### 1.3.4  Causality and Do Calculus in Airbnb Pricing

In addition to regression analysis, causality and do calculus, as developed by Judea Pearl, are powerful tools for analyzing Airbnb pricing models [7]. These tools have also been used to analyze hotel cancellations, as randomized controlled trials can be too costly or unethical [4] By applying causality and do calculus, analysts can determine the factors that have the most significant impact on Airbnb rental prices and develop personalized pricing models that benefit both property owners and travelers. Further research is needed to explore these tools' potential and develop more accurate and precise pricing models for the sharing economy.

## 1.4  Goal and Objectives

### 1.4.1  Research Goal

This project aims to develop a machine-learning model that accurately predicts Airbnb rental prices based on various factors, including property location, size, and amenities. By doing so, we aim to provide an effective tool for property owners to maximize their profits and for renters to make informed decisions on renting properties at fair prices.

### 1.4.2    Associated Objectives

The specific objectives are:

1. Collect Airbnb data from different cities across the United States, including information on property location, size, amenities, and rental prices.

2. Store collected data into structured Pandas DataFrame and separate each DataFrame for each area.

3. Get familiar with preprocessing techniques, machine/deep learning algorithms, libraries, and packages for natural language processing and regression analysis.

4. Preprocess and format the Airbnb data for feeding into the machine learning model.

5. Develop and test multiple regression models to identify the most effective model for predicting Airbnb rental prices.

6. Analyze the importance of different factors in predicting rental prices and provide a quantitative importance score for each variable.

7. Implement and fine-tune the selected regression model to achieve high accuracy and reliability in predicting rental prices.

8. Validate and evaluate the model's accuracy and effectiveness using internal and external testing datasets.

9. By achieving these objectives, we aim to provide a scientifically validated and reliable tool to determine fair and competitive rental prices for the rental market. Additionally, this tool can benefit real estate investors by providing insights into the potential profitability of properties in the market.

## 1.5    Project Organization

In the following chapters, we will provide a detailed overview of our project to predict Airbnb rental prices using machine learning algorithms.

- Chapter 2 will present the requirements analysis, consisting of both functional and non-functional requirements. We will outline the specific features and capabilities that our model must have, as well as any performance, usability, and security requirements.

- Chapter 3 will explain the system architecture and provide detailed designs for each model component. This will include the data collection and preprocessing steps, the feature engineering process, and the regression analysis model. We will also explain how each component interacts with one another to produce the final prediction.

- Based on the previous two chapters, Chapter 4 will provide information on the implementation of the model. We will explain the specific programming languages, libraries, and tools we will use to develop the model and describe the testing and debugging process.

- Chapter 5 will consist of the experimental results of our model, including its accuracy and reliability in predicting Airbnb rental prices. We will provide a detailed analysis of the results, including any limitations or challenges we encountered during the testing process.

- In Chapter 6, we will conclude the project and outline possible extensions or future directions for the research. We will also reflect on the project's overall success and identify any lessons learned or areas for improvement.

Overall, our project aims to provide a scientifically validated and reliable tool for the rental market to determine fair and competitive rental prices. We are confident that our structured approach will lead to a successful outcome.

# CHAPTER 2

# REQUIREMENTS ANALYSIS

This chapter will discuss the functional and non-functional requirements of our Airbnb rent prediction model. The functional requirements will cover the user-defined needs, the behavior of our application, and the necessary components to ensure proper functionality. The non-functional requirements explain the features we will implement to increase usability but are not required for the functionality of our application.

## 2.1    Functional Requirements

1. A dataset that contains more than 9000 Nashville listings of Airbnb properties to develop an understanding of how different variables affect the outcome of an Airbnb Price. It will include Airbnb-listed properties from the Nashville area.

   (a) Each row will contain a property that is listed on Airbnb.

   (b) Each column has attributes regarding property details such as price, number of beds, and customer feedback rating.

2. Structured tabular data from the database will be processed and curated.

   (a) The raw data needs to be prepared for predictive modeling.

   (b) The data must be imputed for missing values.

   (c) Not all columns will be factors to price, such as the host ID or host link, and need to be removed before training the data.

   (d) Columns presenting numeric and categorical variables will be kept and processed for predictive modeling.

3. A regression model to determine the relationships between variables.

   (a) The software shall develop and compare multiple advanced regression models to identify the most accurate model for predicting Airbnb rental prices.

   (b) The software shall learn the importance of different factors in property attributes in predicting rental prices and provide a quantitative relative importance score for each variable.

   (c) The software should provide a range of fair rental prices learning from a set of regression models.

   (d) The software shall provide a causal relationship between the property attributes and rents to provide insights into the cause and effect related to rental price rise and drop.

4. A graphical user app is to be created to allow user input.

   (a) The user should be able to access the software using a simple mouse click and see the application load.

   (b) The user will be able to enter property attributes to obtain a range of fair prices for a property.

   (c) When the user submits the information, the software should show what attribute played a major role in determining the rent.

5. Perform a comparison of various regression and causality tests to determine the most accurate method, presenting the results in a table.

   (a) The software shall evaluate different regression and causality tests, including their respective Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Squared Error (MSE), R-Squared, and Adjusted R-Squared values.

   (b) The comparison results will be presented in a table format, allowing for easy identification of the most accurate method for predicting rental prices.

(c) By identifying the most accurate method, the software ensures the highest quality model is utilized for rental price predictions.

## 2.2     Non-Functional Requirements

1. The software application will be developed and debugged on Jupyter notebook, but the final product will run using an equivalent python script.

   **Rationale:** The Python programming language has one of the largest communities because of the public libraries/packages available relating to Machine Learning algorithms.

2. Ensuring an R2 score above .70 (high)

   **Rationale:** Setting an R2 score above .70 will ensure accurate rental price predictions. This threshold significantly improves over the baseline and reduces pricing prediction errors.

3. The system should run on any desktop environment that runs python.

   **Rationale:** To make the software easily accessible, the application must be able to run in today's most common operating systems.

4. The end user needs to be able to enter property details to get a price result.

   **Rationale:** While it is good to prove the model is successful with a high R2 correlation, we need to make an application that the end user can easily enter in details to give back a prediction.

5. Acceptable confidence interval X +/ 95% sd.

   **Rationale:** The program's requirement to estimate rental prices within an acceptable confidence interval X +/- 95% sd ensures precise and reliable predictions, improving customer satisfaction and enabling property owners to set appropriate rental prices.

6. Regression and Causality Comparison Results

   **Rationale:** Comparing different regression and causality tests helps identify the most accurate method for predicting rental prices, ensuring the highest quality model is utilized.

CHAPTER 3

DESIGN

This chapter outlines our data-driven rental pricing model's system architecture and detailed design. The model will be developed using machine learning algorithms, focusing on data collection, data preprocessing, and model training and testing. In addition, we will be exploring causality by incorporating Microsoft Dowhy library to demonstrate causal discovery into our model.

## 3.1    System Architecture

Our rental price prediction model starts with data gathering and processing. First, we found a data set with around 134,000 Airbnb listings stored in a Pandas DataFrame containing property location, size, amenities, and rental prices.

The next subsystem handles data cleansing and preprocessing. We will remove irrelevant columns, address missing values, and convert categorical variables into numerical values. Feature engineering will create new features from raw data and select the most relevant ones for prediction. The cleaned data is fed into machine learning algorithms, such as linear regression and gradient-boosting tree regression. Hyperparameter tuning will optimize the model's performance.

The causal discovery analysis component will explore relationships between property attributes and rental prices, identifying causal factors influencing rental prices and providing valuable insights. While this method does use linear regression testing in a step, it goes beyond correlation to model the data to understand causal relations.

The final component is the user interface and visualization, offering an accessible platform for users to input property data and receive predictions and insights about rental prices. Visualization techniques like bar charts, scatter plots, and heatmaps will present the results intuitively and appealingly. In addition, as a matter of scientific exploration of knowledge and highlighting the

difference between correlation and causation, the app will break down each prediction model to give the user the best context to give confidence to the results.

In summary, our rental price prediction model comprises interconnected subsystems providing an efficient and accurate solution for predicting Airbnb rental prices. The process includes data gathering and processing, preprocessing, feature engineering, machine learning algorithms, causal discovery analysis, user interface, and visualization components. In addition, using traditional correlation models and the new science of causality, we can demonstrate how data science attempts to move beyond correlation and understand the causality of results to make better decisions.

Figure 3.1: Data processing model.

## 3.2 Detailed Design

The system's design is based on a modular approach, ensuring that each component can be independently developed, tested, and optimized. This structure allows for easy modification and expansion of the system as new data sources and algorithms become available. The design encompasses the following subsections:

### 3.2.1    Data Collection

Data collection involves gathering Airbnb listings from various cities across the United States. The data set is a subset of all Airbnb listing in and is stored as raw data. The data set will include information on property location, size, amenities, and rental prices. We will extract this data from a publicly available website for local storage for analysis.

### 3.2.2    Data Cleansing

In this step, we will preprocess and clean the raw data. This involves imputing missing values, removing irrelevant columns, and converting categorical variables into numerical values. For example, removing non-numerical columns and finding important variables/features in the entire data set. This code will create binary columns for each unique amenity in your 'Amenities' column and merge them back into the original data frame. Each row will have a "1" if the amenity is present and a "0" if not. It will be vital to avoid the curse of dimensionality, which means it will be important to understand feature selection and pre-processing to build an enhanced model [10]. Data cleansing ensures that the machine learning models can effectively process the data and provide accurate predictions.



Figure 3.2: Data cleaning model.

### 3.2.3      Causal Discovery Analysis

Causal discovery analysis allows us to identify the factors that have the most significant impact on Airbnb rental prices by addressing confounding issues. Confounding issues arise when a variable influences the predictor and the outcome, leading to spurious correlations and potentially biased results. By employing the dowhy library, we will apply causality and do-calculus techniques to determine the causal relationships between property attributes and rental prices while accounting for potential confounders.



Figure 3.3: Directed acyclic airbnb causality model.

Addressing confounding issues is crucial for reliable and accurate causal insights; according to Judea Pearl in The Book of Why, he explains that counterfactual analysis will allow scientists to make more precise and definite statements than before [3]. We will identify potential confounding variables affecting property attributes and rental prices to achieve this. First, we can represent the variables' underlying causal structure using existing causal analysis graphical tools. Then, by adjusting the confounding factors, we can estimate the causal effects of the property attributes on rental prices.

### 3.2.4      Linear Regression Testing

Linear regression analysis is used to predict the value of a variable based on the value of another variable. We will develop a linear regression model to analyze the relationship between property

attributes and rental prices. This model will be used as a baseline to compare the performance of other machine learning algorithms. We will evaluate the model's accuracy using metrics such as R2 score and mean squared error.

### 3.2.5 Decision Tree Regression Testing

Decision tree regression is a type of machine learning algorithm that recursively splits the data into subsets based on the most informative features, creating a tree-like structure. In this section, we will implement a decision tree regression model and compare its performance to the linear regression model. We will also fine-tune the model's hyperparameters to achieve the best possible accuracy.

### 3.2.6 Gradient Boosting Tree Regression Testing

Gradient boosting tree regression is a powerful machine learning algorithm that combines multiple weak learners to create a more accurate model. This section will implement a gradient-boosting tree regression model and compare its performance to the linear regression model. In addition, we will fine-tune the model's hyperparameters to achieve the best possible accuracy.

### 3.2.7 HistGradientBoostingRegressor Testing

HistGradientBoostingRegressor is a histogram-based gradient boosting algorithm available in the scikit-learn library. It is designed for large-scale data sets and is an efficient alternative to the traditional gradient boosting method. In this section, we will implement a HistGradientBoostingRegressor model and compare its performance to the linear regression, decision tree, gradient boosting tree, random forest, and XGBoost models. We will also fine-tune the model's hyperparameters to achieve the best possible accuracy.

### 3.2.8 Random Forest Regression Testing

Random forest regression is an ensemble learning method that combines multiple decision trees to improve prediction accuracy and reduce over fitting. In this section, we will implement a random forest regression model and compare its performance to the linear regression, decision tree, and gradient boosting tree models. We will also optimize the model's hyperparameters to obtain the highest accuracy.

### 3.2.9 XGBoost Regression Testing

XGBoost (eXtreme Gradient Boosting) is an optimized gradient boosting library designed for high-performance machine learning tasks. In this section, we will implement an XGBoost regression model and compare its performance to the linear regression, decision tree, gradient boosting tree, and random forest models. Furthermore, we will fine-tune the model's hyperparameters to achieve the best possible accuracy.

CHAPTER 4

IMPLEMENTATION

In this chapter, we delve into the design and implementation of our Airbnb pricing analysis. The process begins with the preprocessing and categorization of Airbnb listing data. This involves handling missing data, categorizing amenities, and converting certain categorical variables into numerical format. Subsequently, the chapter outlines the steps to download and prepare the Nashville Airbnb dataset. Two primary analysis techniques, traditional statistical methods and causal inference models, are employed on this dataset. For a comprehensive analysis, various amenities are grouped, and their potential influence on listing prices is investigated. Lastly, a causal model is used to understand the true impact of these amenities on price, taking into account confounding variables.

## 4.1 Software Environments

Visual Studio Code (VS Code) was the primary tool of choice throughout this project, with its seamless support for Jupyter notebooks and Python scripting. The following specifics were adhered to:

For the majority of the modeling, Python version 3.9.17 was employed. Due to library dependencies, Python version 3.7.16 was utilized for the causal analysis as it necessitates an older version of numpy. Visual Studio Code offers an integrated development environment, facilitating efficient coding, debugging, and testing capabilities. Its versatility in supporting multiple Python versions made it particularly suitable for this project. VS Code's integrated terminal and extensions further streamlined the data analysis and modeling process.

## 4.2 Front-end Development

To ensure that the insights and results from the Airbnb pricing analysis are accessible and interactive, a decision was made to develop a local web application. Flask, a lightweight web application framework written in Python, was chosen for this purpose due to its simplicity, flexibility, and compatibility with our Python-based analytical models.

***Features of the Flask Web Application:***

1. Interactive Dashboard: Users can select specific amenities or features to understand their influence on listing prices.

2. Predictive Tool: Leveraging the models developed, users can input specific listing features to receive a suggested price range.

3. Simplicity: Flask provides a straightforward way to develop web applications, making it perfect for projects that require a simple interface with powerful backend capabilities.

4. Integration with Python: Given that our analysis and models are developed in Python, Flask seamlessly integrates with our existing codebase.

5. Development Speed: Flask's minimalistic approach allows for rapid development and testing, ensuring that the web application can be iteratively improved.

6. Local Deployment: Flask supports easy local deployment, enabling stakeholders to run the web application on their machine without any intricate setup. By integrating Flask, the intention is to provide a user-friendly interface where stakeholders can easily understand the dynamics of Airbnb listing prices in Nashville and make informed decisions accordingly.

## 4.3　Libraries and Tools

In this section, we will go over the general libraries, how we handle the data, visualization tools, and the models that we use for machine learning.

### 4.3.1　General Utilities:

1. warnings: Used to handle and suppress warnings generated during code execution.

2. os: Essential for interacting with the operating system, especially for directory manipulations.

3. re: Provides support for regular expressions, useful for data cleaning and text processing.

4. Counter from collections: Helps in counting occurrences of elements, particularly useful for categorical data analysis.

### 4.3.2　Data Handling and Analysis

1. pandas: Essential for data manipulation and analysis.

2. numpy: Provides support for large arrays and matrices, along with a collection of mathematical functions.

### 4.3.3　Visualization

1. matplotlib: Provides a wide array of tools for generating plots, histograms, and other visual representations.

2. seaborn: Built on top of matplotlib, it offers a higher-level interface for statistical graphics.

### 4.3.4　Modeling and Machine Learning

1. xgboost: An optimized gradient boosting library known for its performance and efficiency.

2. scikit-learn: A comprehensive library containing tools for machine learning and statistical modeling, including:

3. LinearRegression: Regression algorithms.

4. DecisionTreeRegressor, RandomForestRegressor, GradientBoostingRegressor, HistGradient-BoostingRegressor: Tree-based regression models.

   (a) DecisionTreeRegressor, RandomForestRegressor, GradientBoostingRegressor, HistGra-dientBoostingRegressor: Tree-based regression models.

   (b) train_test_split: For splitting datasets into training and testing sets.

   (c) r2_score, mean_squared_error, mean_absolute_error: Metrics for evaluating regression models' performance.

   (d) permutation_importance: For assessing the importance of features in the model.

   (e) KFold: Used for cross-validation.

   (f) dowhy: Provides tools for estimating causal effects using machine learning methods, including:

### 4.3.5    Causal Analysis:

1. dowhy: Enables causal inference by providing a unified interface for model selection, esti-mation, and validation.

2. While we don't use the causal discovery code in the program, there are additional libarys that will be needed to run causal discovery.

### 4.3.6    Others:

1. logging: Useful for generating logs, aiding in debugging and keeping track of the applica-tion's flow.

2. IPython.display: Assists in displaying objects like images directly in Jupyter notebooks or IPython environments.

3. joblib: Provides utilities for saving and loading Python objects, enhancing the reusability of large datasets or trained models.

These libraries and tools are crucial in ensuring that the analysis is robust, the models are accurate, and the application is interactive and user-friendly.

## 4.4 Understanding the Data and Model Building

One of the most crucial stages in a data science project is understanding the nature of the data at hand. This involves diving deep into the dataset, analyzing distributions, identifying trends, and observing relationships among variables. With a dataset downloaded it needs to be imported into our python dataframe.

### 4.4.1 Data Exploration

The dataset for this project, 'nashvilleDF', contains detailed listings of properties available for rent on Airbnb in Nashville. One of the key variables of interest is the 'price', which represents the nightly rate for each property. Before any form of analysis could proceed, the 'price' variable needed to be cleaned and transformed into a numeric format:

### 4.4.2 Analysis by Room Type

A fundamental question was understanding how prices vary by the type of room or property. This can provide insights into what guests are willing to pay more for and how different property types are positioned in the Nashville market.

```
1  # Calculate average price for each room type
2  avg_prices = nashvilleDF.groupby('room_type')['price'].mean()
3
4  # Plot
```

```
 5  plt.figure(figsize=(10,6))
 6  avg_prices.plot(kind='bar', color=['red', 'blue', 'green', 'yellow'])
 7  plt.ylabel('Average Price')
 8  plt.xlabel('Room Type')
 9  plt.title('Average Price by Room Type')
10  plt.xticks(rotation=45)
11  plt.tight_layout()
12  plt.show()
```

Listing 4.1: A bar chart was generated to visualize the average price for each room type.



Figure 4.1: Visualizing the data to understand how different features affect price.

### 4.4.3 The visualization revealed four distinct categories of room types:

1. Entire home/apartment

2. Hotel room

3. Private room

4. Shared room

From the analysis, hotel rooms emerged as the most expensive on average. This observation could be attributed to the premium services and amenities often associated with hotels. The other

room types also had distinct average prices, shedding light on the value perceptions of different accommodation types in the Nashville Airbnb market.

## 4.5 Building the Model

With a clearer understanding of the dataset's structure and the relationships between variables, the next step was to construct predictive models. The goal was to predict the price of a listing based on its features. The journey of model building, from feature selection to model evaluation, is detailed in the subsequent sections.

### 4.5.1 Property Type Transformation

The initial exploration of the 'property type' column, Table 4.1, revealed a myriad of categories, many of which were overlapping or closely related. The goal was to consolidate these into broader, more meaningful categories. The first observation was that numerous property types had the term "Private" in their descriptions.

```
1
2 print(nashvilleDF['property_type'].value_counts())
3
4 To streamline this, all types containing the term "Private" were grouped into a single
      "Private room" category:
5 private_mask = nashvilleDF['property_type'].str.contains('Private')
6 nashvilleDF.loc[private_mask, 'property_type'] = 'Private room'
```

Listing 4.2: Property type grouping.

### 4.5.2 Visualizing the Transformed Data

As shown in Figure 4.2, the 'property type' column was streamlined, and the next step was to visualize how the average price varies across these categories.

Table 4.1: Accommodation types and their counts.

| Accommodation Type | Count |
|---|---|
| Entire home | 3206 |
| Entire rental unit | 1642 |
| Entire condo | 1404 |
| Entire townhouse | 1010 |
| Private room | 723 |
| Entire guest suite | 304 |
| Hotel | 278 |
| Entire guesthouse | 205 |



Figure 4.2: Property type vs cost.

### 4.5.3    Host Response Time

Shown in Table 4.2, one of the critical features in the dataset is the host's response time, which can influence a potential renter's decision.

Table 4.2: Response time by the host.

| Response Time | Count |
|---|---|
| within an hour | 7089 |
| within a few hours | 510 |
| within a day | 191 |
| a few days or more | 39 |

After validating a sizable amount of response time it is important to make this feature more model-friendly, as such the response times were encoded numerically, as shown in Table 4.3.

Table 4.3: Response time by the host.

| Response Time cleaned | Count |
|---|---:|
| within an hour | 1 |
| within a few hours | 2 |
| within a day | 3 |
| a few days or more | 4 |

### 4.5.4    Boolean Columns Conversion

Several columns in the dataset had boolean values represented as 't' for true and 'f' for false. These columns include:

1. host is superhost

2. instant bookable

3. has availability

4. host has profile pic

5. host identity verified

```
nashvilleDF['host_is_superhost'] = [
    1 if x == 't' else 0 for x in nashvilleDF.host_is_superhost
]
```

Listing 4.3: Example of normalizing true false values.

All these columns were converted to a binary representation where 't' was transformed to 1 and 'f' to 0, as shown in Listing 4.3. This transformation aids in the modeling process by providing a numerical representation of boolean data.

### 4.5.5    Handling Bathroom and Bedroom Data

The 'bathrooms text' column in the dataset contains textual information about the number of bathrooms in each property. To create a model-friendly feature, this textual information was converted

to a numerical format, as shown in Listing 4.4. The numeric portion was extracted from the 'bath-rooms text' column, and the resultant values were stored in a new column called 'bathrooms'. The original 'bathrooms text' column was then dropped to avoid redundancy.

```
1 nashvilleDF['bathrooms'] =
      nashvilleDF['bathrooms_text'].str.extract('(\d+\.?\d*)').astype(float)
2 # Drop the 'bathrooms_text' column
3 nashvilleDF = nashvilleDF.drop(columns=['bathrooms_text'])
```

Listing 4.4: Extract numeric portion from the 'bathrooms text' column and assign to 'bathrooms' column.

### 4.5.6    Initial Distribution of Bathrooms

The distribution of the 'bathrooms' column in our dataset displayed a concentration of properties with between 1 and 4 bathrooms. However, certain properties reported an unusually high count, with some even reaching 18 bathrooms. This led to the suspicion of potential outliers in the data.



Figure 4.3: Distribution of bathrooms highlighting potential outliers.

Figure 4.4: Distribution of bedrooms, for comparison.

### 4.5.7    Outlier Detection and Removal

To better visualize and confirm the presence of outliers in the 'bathrooms' distribution, we employed a boxplot. Further, a custom Python function, 'findOutliers', was developed to algorithmically identify and flag outliers. This function operates based on the interquartile range (IQR) method, marking data points outside the range defined by 'q75 + lims calar * iqr' as outliers. Here, 'lim scalar' is a user-defined scalar multiplier, and 'q75' represents the 75th percentile of the data.

```python
def findOutliers(df, column, lim_scalar=4):
    q25, q50, q75 = df[column].quantile(q=[0.25, 0.5, 0.75])
    iqr = q75 - q25
    # Maximum value to be considered as an outlier
    max_ = q75 + lim_scalar * iqr
    # Identify the outliers
    outlier_mask = [True if x > max_ else False for x in df[column]]
    print('{} outliers found out of {} data points, {:.2f}% of the data. Max threshold is
    {}'.format(
        sum(outlier_mask), len(df[column]),
        100 * (sum(outlier_mask) / len(df[column])), max_))
```

```
11      return outlier_mask
```

Listing 4.5: Function to identify outliers based on IQR method.

### 4.5.8    Data Pruning

Specific columns in the dataset were not relevant to the analysis or contained redundant informa-
tion. These columns were thus dropped to streamline the dataset. The columns removed include
URLs, names, descriptions, location-specific information, and other such columns not directly
contributing to the modeling process.

### 4.5.9    Data Cleaning

The price column was crucial for the analysis, and its distribution was inspected. Any listings
with prices outside the range of 200 to 1000 were removed to narrow down the focus to a more
standard set of accommodations. Additionally, some columns contained percentage symbols or
other non-numeric characters. These were cleaned, and missing values in the dataset were imputed
using appropriate strategies, such as median or mean imputation.

```
1 nashvilleDF['price'] = nashvilleDF['price'].str.replace('$', '').str.replace(',',
      '').astype(float)
```

Listing 4.6: Example of cleaning up features.

### 4.5.10    Visualization of Price Distribution

After the cleaning and pruning, the distribution of the price column was visualized. The resulting
histogram showed a bell-like curve, shown in Figure 4.5, indicating a standard distribution of prices
in the dataset. This visualization can provide insights into the most common price range for Airbnb
listings in the area.

Figure 4.5: Distribution of price of all the properties.

### 4.5.11    Model Building and Feature Importance

The last step of our analysis involved building a model to predict Airbnb prices based on the cleaned and processed features. For this task, we chose the HistGradientBoostingRegressor, an ensemble method known for its ability to capture complex relationships and handle large datasets efficiently.

1. Cross-Validation: To ensure the robustness of our model, a 10-fold cross-validation was implemented. This approach divides the dataset into ten parts: in each iteration, nine parts are used for training, and one part is used for testing. This iterative process helps in understanding the model's average performance over different parts of the dataset and reduces the chances of overfitting to a particular subset.

2. Feature Importance: One of the major challenges in building machine learning models is understanding which features play a significant role in predictions. Using the permutation importance function, we can rank features based on their importance. This function works by randomly shuffling a feature's values and measuring how much the model's performance drops. If shuffling a feature's values causes a significant drop in performance, that feature is

deemed important.

3. Result:

The model was trained and evaluated on each fold of the data. For every iteration:

The model's accuracy on the test set was printed.

The features were ranked by their importance.

The average accuracy across all folds provided an aggregate measure of the model's performance.

Additionally, the standard deviation of accuracy scores gave an indication of the model's consistency across different parts of the dataset.

To provide further clarity on the significance of different features, an average rank of the top features across all folds was computed. This average rank highlighted the most consistently important features across all iterations.

### 4.5.12    Implications for User Interface

The insights from the feature importance analysis are invaluable for the development of the user interface. By understanding which features significantly impact price, the UI can be tailored to emphasize these key features. This ensures that users are prompted to provide the most relevant information when seeking price predictions. For instance, if the number of bathrooms consistently ranks as a top feature, it would be wise to prominently display this option in the predictive tool of the web application. On the other hand, less impactful features could be relegated to advanced settings or optional inputs. We chose to build a web application using HTML, CSS, and Python Flask. All three of these languages work together to create a functioning application.

### HyperText Markup Language

HTML, or HypertText Markup Language, was used to build the structure and presentation of the web application, including text, images, links, forms, and multimedia. HTML uses tags and attributes to describe the elements and their relationships within a web page.

```html
<div class="topnav">
        <a class="active" href="{{ url_for('predictor') }}">Predictor</a>
        <a href="{{ url_for('about') }}">About</a>
</div>
```

Listing 4.7: Web pade tool bar code.

The tag "form" allows users to input and submit data to a web server for processing. HTML forms serve as a means for the user to interact with the web application, provide information, and trigger actions. Listing 4.8 demonstrates what these forms tags look like in code.

```html
<form method="POST" action="{{ url_for('predictor') }}">
        {{ form.hidden_tag() }}
        <div class="row">
            <div class="column">
                {{ form.accommodates.label }}<br>{{
    form.accommodates(class="input-field") }}
                {{ form.bathrooms.label }}<br>{{ form.bathrooms(class="input-field") }}
                {{ form.nights_staying.label }}<br>{{
    form.nights_staying(class="input-field") }}
                {{ form.bedrooms.label }}<br>{{ form.bedrooms(class="input-field") }}
                {{ form.fireplace.label }}<br>{{ form.fireplace() }}
                {{ form.hot_tub.label }}<br>{{ form.hot_tub() }}
            </div>
            <div class="column">
                {{ form.cable.label }}<br>{{ form.cable() }}
                {{ form.wifi.label }}<br>{{ form.wifi() }}
                {{ form.review_scores_value.label }}<br>{{
    form.review_scores_value(class="input-field") }}
                {{ form.property_type.label }}<br>{{
    form.property_type(class="input-field") }}
                {{ form.room_type_num.label }}<br>{{
    form.room_type_num(class="input-field") }}
```

30

```
18                {{ form.room_type_num.label }}<br>{{
    form.room_type_num(class="input-field") }}
19          </div>
20       </div>
21       {{ form.submit(class="submit-button") }}
22    </form>
```

Listing 4.8: HTML forms code.

Every page of the web page has its own unique HTML file. HTML is ofter combined with CSS to create interactive and visually appealing web pages.

### *Cascading Style Sheet*

CSS, or Cascading Style Sheets, was used to control the visual presentation and layout of HTML elements on the web page. CSS allows the web application to adapt to various screen sizes and devices, customize the look and feel of the user interface, and create animations and interactions without relying on Javascript. As shown in Listing 4.9, CSS styles can be applied directly within a HTML document or linked externally to it, making it easier to maintain and update the appearance of a website.

```
1    <link rel="stylesheet" type="text/css" href="{{ url_for('static',
    filename='style.css') }}">
```

Listing 4.9: Implementing the CSS style sheet to the HTML web page.

CSS allowed us to define how text, images, links, and elements should look on the website. CSS works by selecting HTML elements and applying style to them using selectors and properties. Selectors target specific elements, such as headings or paragraphs, and properties define the style, like the font size or background color.

```
1 body {
2     height: 100%;
3     background-color: #51504c;
4     color: white;
5     font-family:Cambria, Cochin, Georgia, Times, 'Times New Roman', serif;
```

31

```
 6 }
 7
 8 input[type=text] {
 9     width: 50%;
10     padding: 12px 20px;
11     margin: 8px 0;
12     box-sizing: border-box;
13     border: 2px solid black;
14     border-radius: 4px;
15     background-color: rgba(128, 128, 128, 0.281);
16 }
17
18 button {
19     background-color: #6e6e6e51;
20     border: 2px solid black;
21     color: rgb(250, 248, 248);
22     padding: 15px 32px;
23     text-align: center;
24     text-decoration: none;
25     display: inline-block;
26     font-size: 12px;
27 }
```

Listing 4.10: CSS code.

### Python Flask

Python Flask was used as the web framework for our web application. Flask is a lightweight and versatile web framework for building web applications in Python. Flask is often referred to as a micro-framework because it provides the essential components for building web applications without imposing a rigid structure or including a lot of built-in functionality. Flask allows you to define URL routes, which determine how different URLs or endpoints of the web application should be handled.

```
1 @app.route('/price')
2 def price():
3     estimated_price = request.args.get('estimated_price', 'N/A')
4     mae = request.args.get('mae', 'N/A')
```

```
5     return render_template('price.html', estimated_price=estimated_price, mae=mae)
```

Listing 4.11: Route example code.

Flask supports HTTP methods like GET, POST, PUT, DELETE, etc., which enables you to handle different types of requests from web clients.

```
1 @app.route('/predictor', methods=['GET', 'POST'])
```

Listing 4.12: HTTP methods code

Flask includes a templating engine, called Jinja2, that allows you to create dynamic HTML templates. These templates can be populated with data from your Python code, enabling you to generate dynamic web content.

# CHAPTER 5

## RESULTS

The foundation of any robust data science project lies in the quality of data, the depth of its analysis, and the insights gleaned from it. This chapter dives deep into the results obtained from our rigorous exploration of a Nashville Airbnb dataset. By utilizing traditional statistical techniques and cutting-edge machine learning models, we are prepared to share inherent patterns, namely relationships between input features and predictive outcomes, that can provide actionable insights for stakeholders. Here, we will walk through the steps of our analytical journey, from identifying the model to creating a custom model for the problem, primarily using regression models to predict fair market pricing. The model-driven outcomes are further interpreted using state-of-the-art causal inferencing tools to rule out spurious correlation in data and provide more transparent and trustworthy predictors of model outcomes.

## 5.1    Regression analysis and results

To begin a machine learning project, one must first explore the data to decide what model to choose and build for the problem in hand. With a large variety of models available, the general practice of grooming the data is standard and allows us to test six different regression models for our project quickly. These regression models provide a function that describes the relationship between one or more independent variables and a dependent target variable. We tested various advanced regression models, including linear regressor, decision tree, Gradient Boost- ing tree, and HistGradientBoosting regressor. We chose the best regression model based on the testing results of the highest R2 score of 72.22. After testing, we decided that the HistGradientBoosting regressor was the best regression model for our project.

Figure 5.1: The results of the Regression Model testing.

## 5.2    Feature Selection Results

The data dimensionality or the number of features is critical for model performance. While more features can provide more information about the relationship, it can lead to 'the curse of dimensionality' effect, hurting the model performance. Concurrently, we also wanted to assess the computational efficiency of the model by observing the time taken to process various numbers of features. We knew that every additional feature would take additional computation time; with that in mind, we want to see when the results start to level out. Therefore, identifying the most important combination of features, namely feature selection, is an imperative step in machine learning modeling.

While our goal was to narrow down to approximately 20 distinct features, certain features, like 'property type', required one-hot encoding. This encoding process meant that to differentiate

Figure 5.2: We can see the importance of the first 20 features and then see a more even increase in its feature importance but with greater time.



Figure 5.3: Knowing that each feature adds a linear time increase but adds no value to the r2 score allows us to reduce the features, resulting in faster computational time for the user using our software.

among categories such as a house, an apartment, or a motorhome, we had to introduce multiple feature columns. As a result, the actual number of feature columns utilized was 28, but they represent 20 unique features. When we re-ran the model, linear regression showed unexpectedly

negative results. However, with the HistGradientBoosting model, we still achieved our target R-squared score of 70



Figure 5.4: Updated results of the Regression Model testing.

Table 5.1: Final r2, MAE, MSE, and RMSE score.

| Model | MAE | MSE | R2 | RMSE |
|---|---|---|---|---|
| DecisionTree | 83.962371 | 19214.576281 | 0.285639 | 138.616652 |
| GradientBoosting | 68.014278 | 9572.757100 | 0.644103 | 97.840468 |
| HistGradientBoosting | 60.517344 | 7977.994567 | 0.703393 | 89.319620 |
| Linear | 123.133481 | 26900.691076 | -0.000116 | 164.014301 |
| RandomForest | 58.805755 | 8294.641718 | 0.691621 | 91.074924 |
| XGBoost | 61.770700 | 8923.464897 | 0.668243 | 94.464093 |

## 5.3 Causality

In the progression of our study, a central ambition was to integrate causal inference, allowing us to delve deeper into the underlying mechanisms of our data. By posing a causal question, we sought not just to establish correlations but to determine the causes and effects that drive the observed relationships. After formulating our causal query, we proceeded to generate an estimate that quantifies the magnitude and direction of this causal relationship. However, an estimate alone is not sufficiently robust. To ensure the veracity of our findings, we introduced a refutation test, producing a refutation score. A statistically significant refutation score acts as a litmus test for our model's accuracy. It indicates that our model's predictions are not merely the outcome of chance or spurious correlations but are likely grounded in genuine causal relationships. As we delved deeper into the importance of each feature, it paved the way for a clearer understanding of the potential causal questions that could be posed, thereby enhancing the depth and rigor of our analysis. With this foundational understanding, we then sought to pinpoint specific causal relationships that might have a tangible impact on property pricing values.

### 5.3.1 Causal Question

Every property has its own story to tell and price that should reflect it, but there are simple things a homeowner could do, like swapping a couch to a sofa bed to increase the number of beds and accommodate in a property. Our causal question was, "If an owner increases the total accommodates, how much more can they rent out their property on average per day in the Nashville market?"

### 5.3.2 Causal Discovery

Domain knowledge is pivotal in constructing an accurate Direct Acyclic Graph(DAG), a foundational component of Judea Pearl's causal model framework. However, since the commencement of our project, several algorithms and models, including Linear Non-Gaussian Acyclic Models (LiNGAM), Greedy Equivalence Search (GES), and the Peter-Clark Algorithm, have emerged.

These methodologies leverage data to predict causal relationships. While these models are progressively becoming more precise, there remains a discernible room for improvement.



Figure 5.5: Subset of LiNGAM Model Visualization: The model identifies a pathway from 'beds' to 'accommodates,' validating one causal inquiry, but fails to connect 'price' to any feature resulting in an incomplete graph.



Figure 5.6: The PC method visualization adeptly clusters features such as 'resort access,' 'pool,' 'hot tub,' and 'grill,' along with 'bathrooms,' 'bedrooms,' 'beds,' and 'accommodation.' However it does not connect all features

Figure 5.7: GES Graph: From our domain perspective, this appears to be the most accurate representation yet, though it's important to note that repeated executions yield varied graphs.

*Summary of Causal Discovery Algorithms*

## LiNGAM Model

The LiNGAM algorithm successfully identifies expected causal directions, such as from 'beds' to 'accommodates', demonstrating its capability in uncovering linear causal relationships. Its failure to establish a connection from 'price' to other variables signals a need for model refinement, possibly requiring additional data or parameter adjustment. The numerical weights on the LiNGAM graph edges are indicative of the magnitude and direction of effects, offering valuable insights into how changes in one variable affect another.

**PC Algorithm**

The PC algorithm reveals potential causal structures in a binary manner, highlighting the presence of causal links without estimating their strengths. This necessitates further analysis, underpinned by domain knowledge, to unravel the underlying causal mechanisms and quantify the strengths of the relationships suggested by the algorithm.

**GES Algorithm**

The GES algorithm generates graphs that align with the best structural fit to the data, according to the chosen scoring function. Despite providing a structural overview, the variability across different executions and the absence of quantitative edge weights underscore the importance of domain expertise. Such knowledge is essential for validating the causal relationships and for enriching the algorithmic findings with contextual depth and insight.

**In summary**, the LiNGAM algorithm is adept at quantifying causal effects, while the PC and GES algorithms excel at outlining the structure of causal networks. Domain expertise plays a pivotal role across all methods for interpreting, validating, and refining the models. An integrated approach, combining multiple algorithmic outputs and informed by expert knowledge, is advocated to develop a robust and nuanced causal understanding.

***Large Language Models Meets Causal DAG***

The intersection of Large Language Models (LLMs) and Causal Discovery is an emerging area of interest among data scientists. A compelling illustration of this synergy can be found in an article penned by Aleksander Molak on Medium.com. Molak, also the author of the book "Causal Inference and Discovery in Python," showcased how integrating ChatGPT with machine learning causal discovery, specifically the PC method can foster a more robust model, courtesy of the LLM's capabilities.[6] We emulated this approach, employing a subset of our features for experimentation. Here, we present our findings contrasting models with and without the influence of an LLM.

Figure 5.8: Metrics without LLM vs. True DAG: On the left, we have the DAG derived without the influence of an LLM, and on the right, the true DAG for reference.



Figure 5.9: Metrics without LLM vs. True DAG: On the left, we have the DAG derived without the influence of an LLM, and on the right, the true DAG for reference.

### Domain Knowledge DAG

Relying on our domain expertise, we crafted a DAG tailored to our dataset. Our next steps involve generating an estimate and a refutation score for our target variable, 'accommodates,' with 'price' as the outcome variable. A salient observation during our experimentation with neural network models was the mutable nature of our graphs upon each iteration. This reinforces our belief that while modern methods offer valuable insights into data relationships, the infusion of domain knowledge remains indispensable.

Figure 5.10: DAG Based on Domain Knowledge: Subset of our graph showing our interpretation and representation of the data, rooted in our expertise.

### 5.3.3    Causal Inference Estimator

After establishing the causal graph and understanding the potential relationships among the features, we moved to the next crucial step in our causality analysis: causal inference. This step enabled us to derive insights regarding the effect of our treatment variable, 'accommodates,' on the outcome variable, 'price.' We utilized a backdoor linear regression method, targeting the average treatment effect (ATE). Here is the Python code we employed: The causal model's results showed a mean value of approximately 21.92 for the effect of 'accommodates' on 'price.' In essence, this suggests that for every unit increase in 'accommodates,' the price could potentially increase by approximately 21.92 dollars after accounting for other variables. This effect was consistent across various conditions and feature combinations, as evidenced by the numerous conditional estimates provided. However, an essential aspect of causal inference is not just estimating the causal effect but validating its authenticity. Therefore, we introduced refutation methods to challenge our initial findings and assess their robustness.

### 5.3.4    Causal Refutation Test

In the realm of causal inference, the reliability of an estimated causal effect is paramount and the least understood step in causal inference. Refutation tests serve as a mechanism to gauge the robustness of these estimates by challenging them under controlled perturbations. For our analysis, we employed two such methods: the 'data subset refuter' and the 'random common cause.'

*Data Subset Refuter*

The Data Subset Refuter method operates by randomly excising a subset of the data and reevaluating the causal effect. If the recalculated causal effect remains congruent with the original, even after the perturbation, it lends further credence to the robustness of our initial estimate. Upon application, the refutation test yielded a new effect of approximately 0.4924 post-data subset removal. This was remarkably consistent with our original estimate of 0.4929 derived from the standard scaled data. With a p-value of 0.96, the minimal difference between the two estimates was statistically negligible. After reverse-engineering the standardization, the real effect translated to an increase of approximately 21.94 for every additional accommodation unit, a slight rise from the original $21.92.

*Random Common Cause*

The Random Common Cause refutation method introduces a random common cause to the dataset and reevaluates the causal effect. The rationale is that if introducing a spurious variable doesn't significantly alter the estimate, then the original effect is likely robust. In our scenario, the effect derived after accommodating the random common cause was The original beta value, represented as "beta original," is equal to 0.4928671844013451 with its .98 score, we know the results are valid.

### 5.3.5    Refutation Test Conclusion

The congruence between our original estimates and those obtained post-refutation affirms the reliability of our causal effect. The minute deviations observed, especially when contextualized with the real-world implications, underscore the robustness of our findings. Such consistency across different refutation techniques bolsters our confidence in the derived insights, emphasizing their potential applicability in real-world scenarios.

```
1  Estimated effect:0.4928598248111045
2  New effect:0.4928671844013451
3  p value:0.98
```

Listing 5.1: Results from Random Common Cause refutation test.

### 5.3.6    Causality Results

Our causal analysis journey, from formulating the right questions to employing modern techniques for causal discovery and validation, underscores the intricate balance between data-driven methods and domain expertise. While algorithms and models provide a structured pathway to unearth causal relationships, the significance of domain knowledge remains unparalleled, acting as a compass in navigating the complex terrains of causality.

Our findings shed light on the potential strategies homeowners in the Nashville market might employ to enhance their property's value. However, as with all models, it's essential to approach the results with a degree of caution and an understanding of the model's limitations. The ever-evolving domain of causal inference promises even more refined tools in the future, paving the way for more nuanced insights and strategies.

## 5.4    Web Application

The web application was designed to provide a user-friendly interface for the purpose of predicting feasible Airbnb rental pricing. The web app was built using the Python Flask framework, with the

front-end built using HTML and CSS. The user will enter features of the property they would like to get a predicted price for using the opening page, shown in Figure 5.11.



Figure 5.11: The screen where the user will enter property features to predict price.

After the user submits the features, their features are put through our prediction model. The model then returns three results, shown in Figure 5.12. It tells the user the predicted price, the mean absolute error, and five similar Airbnb properties. The mean absolute error is the amount the predicted price can vary.

Predictor    Price

Predicted Price: $89.37

MAE: $59.9 (This indicates the average deviation of the predicted price from actual prices. A lower value means the prediction is more accurate.)

## Recommended Airbnb Listings:

Property ID: 917534000000000000
Property ID: 716749000000000000
Property ID: 721723000000000000
Property ID: 720426000000000000
Property ID: 19033669

Figure 5.12: This screen is presented when the user submits their property features.

# CHAPTER 6

# CONCLUSIONS

In this comprehensive analysis of the Nashville Airbnb market, we employed a series of methodological approaches, each underpinned by robust scientific principles. Beginning with the raw dataset, meticulous preprocessing techniques ensured the data's integrity and relevance, laying the groundwork for sophisticated modeling.

Our modeling phase was characterized by rigorous empirical testing across multiple regression models. The HistGradientBoosting Regressor emerged as the optimal model, demonstrating superior performance metrics in capturing the complexities of the dataset. From here we were able to do feature reduction to allow our Python flask app to run faster with minimal, and acceptable,l statistical loss.

Transitioning to causality, we ventured into the intricate domain of causal inference, systematically exploring the cause-effect relationships inherent in our data. Leveraging both state-of-the-art algorithms and established domain knowledge, our analysis identified 'accommodates' as a statistically significant determinant of property prices.

The confluence of advanced analytical techniques, rigorous empirical methodologies, and domain expertise culminated in a set of actionable insights. These insights, grounded in scientific rigor, provide a robust framework for stakeholders to make informed decisions in the Nashville Airbnb market. This study underscores the importance and potential of data-driven, scientifically rigorous approaches in understanding and optimizing complex market dynamics.

# Appendices

# APPENDIX 1

# CREATING THE MODEL

## A.1 Model Creation

```python
1  # Import necessary libraries
2  import warnings
3  import pandas as pd
4  import numpy as np
5  import matplotlib.pyplot as plt
6  import seaborn as sns
7  import xgboost as xgb
8  from sklearn.linear_model import LinearRegression, SGDRegressor
9  from sklearn.tree import DecisionTreeRegressor
10 from sklearn.ensemble import GradientBoostingRegressor, RandomForestRegressor, HistGradientBoostingRegressor
11 from sklearn.model_selection import train_test_split
12 from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
13 import dowhy
14 from dowhy import CausalModel
15 import re
16 from collections import Counter
17
18
19
20 def load_data(filepath):
21     df = pd.read_csv(filepath, header=None, skiprows=0)
22     df.columns = df.iloc[0]
23     df = df[1:]
24     return df
25
26
27
28 # Load and preprocess data
29 file_path = r"C:\Users\nstep\TSU\SeniorProject\DataSet\listings.csv"
30 nashvilleDF = load_data(file_path)
31
32 #nashvilleDf to lower case
33
34 nashvilleDF.head()
35
36 pd.options.display.max_columns = 200
37 # Colors sourced from here: https://usbrandcolors.com/airbnb-colors/
38 bnb_red = '#FF5A5F'
39 bnb_blue = '#00A699'
40 bnb_orange = '#FC642D'
41 bnb_lgrey = '#767676'
42 bnb_dgrey = '#484848'
```

```
43  bnb_maroon = '#92174D'
44  # Create diverging colormap for heatmaps
45  bnb_cmap = sns.diverging_palette(210,
46                                    13,
47                                    s=81,
48                                    l=61,
49                                    sep=3,
50                                    n=16,
51                                    as_cmap=True)
52
53  # Test colors
54  sns.palplot(sns.diverging_palette(210, 13, s=81, l=61, sep=3, n=16))
55
56
57
58  # Create color palette
59  bnb_palette = sns.color_palette(
60      ["#FF5A5F", "#007989", "#8CE071", "#FC642D", "#92174D", "#01D1C1"])
61
62  # Test colors
63  sns.palplot(bnb_palette)
64
65
66  nashvilleDF.shape
67  #nashvilleDF.info()
68
69  #write nashvilleDF.info to csv
70  info = nashvilleDF.info()
71
72
73
74  # Visualize price table, changing them to floats and replacing the commas with a blank
75  prices = nashvilleDF['price'].apply(lambda s: float(s[1:].replace(',','')))
76
77  # Drop listings with a price of zero
78  prices = prices[prices!=0]
79
80  # Log prices
81  log_prices = np.log(prices)
82
83  print(log_prices.describe())
84
85  def plot_hist(n, titles, ranges):
86      """
87      Quick helper function to plot histograms
88      """
89      fig, ax = plt.subplots(n, figsize = (8, 7.5))
90      for i in range(n):
91          d, bins, patches = ax[i].hist(ranges[i], 50, density = 1, color='red', alpha = 0.85)
92          ax[i].set_title(titles[i])
93          ax[i].set_xlabel("Daily Listing Price in Dollars")
94          ax[i].set_ylabel("Frequency")
95      plt.tight_layout()
96      plt.show()
97
98  # Plot histograms of price distribution
```

```
 99  plot_hist(4, ['Distribution of Listing Prices: All Data', 'Distribution of Listing Prices: $0 - $1000',
100                  'Distribution of Listing Prices: $0 - $350','Log Transformed Distribution of Listing Prices: All Data'],
101            [prices, prices[prices <= 1000], prices[prices < 350],log_prices])
102
103
104
105  # Convert relevant columns to datetime format
106  nashvilleDF.last_scraped = pd.to_datetime(
107      nashvilleDF.last_scraped)
108
109  nashvilleDF.host_since = pd.to_datetime(
110      nashvilleDF.host_since)
111
112  nashvilleDF.calendar_last_scraped = pd.to_datetime(
113      nashvilleDF.calendar_last_scraped)
114
115  nashvilleDF.first_review = pd.to_datetime(
116      nashvilleDF.first_review)
117
118  nashvilleDF.last_review = pd.to_datetime(
119      nashvilleDF.last_review)
120
121
122
123  nashvilleDF.host_response_rate = nashvilleDF[
124      'host_response_rate'].apply(lambda s: float(str(s).replace('%', '')))
125
126
127
128  nashvilleDF.property_type.value_counts()
129
130  sns.countplot(y=nashvilleDF['property_type'])
131  plt.ylabel('Property Types')
132  plt.xlabel('Number of Listings')
133  plt.title('Number of Listings by Propery Type')
134
135  # Step 1: Create a mask to identify rows with 'Private' in the 'property_type' column
136  private_mask = nashvilleDF['property_type'].str.contains('Private')
137
138  # Step 2: Update the values in the 'property_type' column based on the mask
139  nashvilleDF.loc[private_mask, 'property_type'] = 'Private room'
140
141  nashvilleDF.replace('Aparthotel','Hotel',inplace=True)
142  nashvilleDF.replace('Room in aparthotel','Hotel',inplace=True)
143  nashvilleDF.replace('Room in hotel','Hotel',inplace=True)
144  nashvilleDF.replace('Room in boutique hotel','Hotel',inplace=True)
145  nashvilleDF.replace('Entire serviced apartment','Entire condo',inplace=True)
146  nashvilleDF.replace('Entire loft','Entire condo',inplace=True)
147
148  private_mask = nashvilleDF['property_type'].str.contains('Shared')
149
150  nashvilleDF.loc[private_mask, 'property_type'] = 'Shared room'
151
152  # Display the updated value counts
153  updated_value_counts = nashvilleDF['property_type'].value_counts()
154  updated_value_counts
```

```
155
156   # Display the updated value counts
157   updated_value_counts_after_drop = nashvilleDF['property_type'].value_counts()
158   print(updated_value_counts_after_drop)
159
160
161   # List of property types to drop
162   drop_list = [
163       "Entire bungalow", "Entire cottage", "Tiny home", "Camper/RV", "Entire cabin",
164       "Entire vacation home", "Entire place", "Shared room", "Entire villa",
165       "Farm stay", "Entire home/apt", "Earthen home", "Barn", "Entire chalet",
166       "Bus", "Shipping container", "Boat", "Tent"
167   ]
168
169   # Drop rows based on the property_type values in drop_list
170   nashvilleDF = nashvilleDF[~nashvilleDF['property_type'].isin(drop_list)]
171
172   # Display the updated value counts
173   updated_value_counts_after_drop = nashvilleDF['property_type'].value_counts()
174   print(updated_value_counts_after_drop)
175
176
177   updated_value_counts_after_drop = nashvilleDF['host_response_time'].value_counts()
178   print(updated_value_counts_after_drop)
179
180   # Host Response Time
181   # We will treat nan values as
182
183   nashvilleDF[
184       'host_response_time'] = nashvilleDF.host_response_time.map({
185           'within an hour':
186           1,
187           'within a few hours':
188           2,
189           'within a day':
190           3,
191           'a few days or more':
192           4,
193           np.nan:
194           5
195       })
196
197
198
199
200   updated_value_counts_after_drop = nashvilleDF['host_response_time'].value_counts()
201   print(updated_value_counts_after_drop)
202
203
204
205   # Display the updated value counts
206   updated_value_counts_after_drop = nashvilleDF['host_is_superhost'].value_counts()
207   print(updated_value_counts_after_drop)
208
209
210   nashvilleDF['host_is_superhost'] = [
```

```
211     1 if x == 't' else 0 for x in nashvilleDF.host_is_superhost
212 ]
213 # Display the updated value counts
214 updated_value_counts_after_drop = nashvilleDF['host_is_superhost'].value_counts()
215 print(updated_value_counts_after_drop)
216
217
218 # Display the updated value counts
219 updated_value_counts_after_drop = nashvilleDF['instant_bookable'].value_counts()
220 print(updated_value_counts_after_drop)
221
222 nashvilleDF['instant_bookable'] = [
223     1 if x == 't' else 0 for x in nashvilleDF.instant_bookable
224 ]
225
226 # Display the updated value counts
227 updated_value_counts_after_drop = nashvilleDF['instant_bookable'].value_counts()
228 print(updated_value_counts_after_drop)
229
230
231 # Display the updated value counts
232 updated_value_counts_after_drop = nashvilleDF['has_availability'].value_counts()
233 print(updated_value_counts_after_drop)
234
235 nashvilleDF['has_availability'] = [
236     1 if x == 't' else 0 for x in nashvilleDF.has_availability
237 ]
238
239 # Display the updated value counts
240 updated_value_counts_after_drop = nashvilleDF['has_availability'].value_counts()
241 print(updated_value_counts_after_drop)
242
243 # Display the updated value counts
244 updated_value_counts_after_drop = nashvilleDF['host_has_profile_pic'].value_counts()
245 print(updated_value_counts_after_drop)
246
247 nashvilleDF['host_has_profile_pic'] = [
248     1 if x == 't' else 0 for x in nashvilleDF.host_has_profile_pic
249 ]
250
251 # Display the updated value counts
252 updated_value_counts_after_drop = nashvilleDF['host_has_profile_pic'].value_counts()
253 print(updated_value_counts_after_drop)
254
255
256 # Display the updated value counts
257 updated_value_counts_after_drop = nashvilleDF['host_identity_verified'].value_counts()
258 print(updated_value_counts_after_drop)
259
260
261 nashvilleDF['host_identity_verified'] = [
262     1 if x == 't' else 0 for x in nashvilleDF.host_identity_verified
263 ]
264
265 # Display the updated value counts
266 updated_value_counts_after_drop = nashvilleDF['host_identity_verified'].value_counts()
```

```python
267 print ( updated_value_counts_after_drop )
268
269
270 # Property Type
271 # One hot encode
272 prop_type_dummies = pd.get_dummies ( nashvilleDF.property_type ,
273                                     prefix='prop' )
274
275 # Merge with df
276 nashvilleDF = nashvilleDF.merge ( prop_type_dummies ,
277                                  left_index=True ,
278                                  right_index=True )
279
280
281 nashvilleDF.property_type.value_counts ()
282
283
284 nashvilleDF.room_type.value_counts ()
285
286
287 # Room Type
288 # Create numerical column for room type
289 nashvilleDF['room_type_num'] = nashvilleDF.room_type.map ({
290     'Entire home/apt':
291     3,
292     'Private room':
293     2,
294     'Hotel room':
295     1,
296 })
297
298 # One hot encode
299 room_type_dummies = pd.get_dummies ( nashvilleDF.room_type ,
300                                     prefix='room' )
301
302 # Drop "hotel room" column as base case
303 del room_type_dummies['room_Hotel room']
304
305 # Merge with df
306 nashvilleDF = nashvilleDF.merge ( room_type_dummies ,
307                                          left_index=True ,
308                                          right_index=True )
309
310
311 nashvilleDF.room_type_num.value_counts ()
312
313
314 # Combine 'Private room' and 'Hotel room' into 'Private/Hotel room'
315 nashvilleDF['room_type'] = nashvilleDF['room_type'].replace(['Private room', 'Hotel room'], 'Private/Hotel room')
316
317 # Display the updated value counts for room_type column
318 updated_value_counts_room_type = nashvilleDF['room_type'].value_counts ()
319 print ( updated_value_counts_room_type )
320
321 # Map the updated room types to numeric values
322 nashvilleDF['room_type_num'] = nashvilleDF.room_type.map ({
```

```
323     'Entire home/apt': 2,
324     'Private/Hotel room': 1
325 })
326
327 # Display the updated value counts for room_type_num column
328 updated_value_counts_room_type_num = nashvilleDF['room_type_num'].value_counts()
329 updated_value_counts_room_type_num
330
331
332 nashvilleDF.head()
333
334 nashvilleDF[nashvilleDF['bedrooms'].isnull()]
335
336
337 # Extract numeric portion from the 'bathrooms_text' column and assign to 'bathrooms' column
338 nashvilleDF['bathrooms'] = nashvilleDF['bathrooms_text'].str.extract('(\d+\.?\d*)').astype(float)
339
340 # Drop the 'bathrooms_text' column
341 nashvilleDF = nashvilleDF.drop(columns=['bathrooms_text'])
342
343
344 #show nashvilleDF['bedrooms'] that are empty
345 nashvilleDF.bathrooms.value_counts()
346
347
348 nashvilleDF.neighbourhood.value_counts()
349
350
351 nashvilleDF.neighbourhood.fillna('Empty',inplace=True)
352
353 nashvilleDF.neighbourhood.value_counts()
354
355
356 #drop neighbourhood from nashvilleDF
357 nashvilleDF.drop('neighbourhood',axis=1,inplace=True)
358 nashvilleDF.neighbourhood_cleansed.value_counts()
359
360 # Extract the district number from the 'neighbourhood_cleansed' column and assign to 'neighbourhood_cleansed_num'
361 nashvilleDF['neighbourhood_cleansed_num'] = nashvilleDF['neighbourhood_cleansed'].str.extract('(\d+)').astype(int)
362
363 # Display the updated value counts for the 'neighbourhood_cleansed_num' column
364 updated_neighbourhood_cleansed_num_value_counts = nashvilleDF['neighbourhood_cleansed_num'].value_counts()
365 updated_neighbourhood_cleansed_num_value_counts
366
367
368 import re
369 from collections import Counter
370
371 # Format amenities column for analysis
372 nashvilleDF.amenities = nashvilleDF.amenities.apply(
373     lambda x: [i.strip() for i in re.sub('[^a-zA-Z,\/\s\d-]*', '', x.lower()).split(sep=',')] if isinstance(x, str) else
       x)
374
375 # Create a flat list of all amenities entries
376 amenities_list = [item for sublist in nashvilleDF.amenities for item in sublist if isinstance(sublist, list)]
377
```

```python
378  # Count amenities occurrences
379  amenity_counts = Counter(amenities_list).most_common()
380
381  # Examine the top 20 amenities
382  top_amenity_counts = amenity_counts[0:20]
383
384  # Look at the 90 least common amenities
385  lowest_amenity_counts = amenity_counts[-90:]
386
387  # Total unique amenities
388  total_unique_amenities = len(set(amenities_list))
389
390  lowest_amenity_counts, total_unique_amenities
391
392
393
394  # Make a list of amenities of interest
395  amenities_of_interest = [x[0] for x in amenity_counts[0:70]]
396
397  #print amenity_counts 1 per row
398  for amenity in amenity_counts:
399      print(amenity)
400
401  print(amenity_counts)
402
403
404  combine_keywords = {
405      "coffee": "coffee",
406      "shampoo": "shampoo",
407      "toaster": "toaster",
408      "crib": "crib",
409      "hot tub": "hot tub",
410      "refrigerator": "refrigerator",
411      "gym": "gym",
412      "resort access": "resort access",
413      "microwave": "microwave",
414      "kitchen": "kitchen",
415      "camera": "camera",
416      "cable": "cable",
417      "grill": "grill",
418      "stove": "stove",
419      "backyard": "backyard",
420      "bluetooth": "bluetooth",
421      "wifi": "wifi",
422      "oven": "oven",
423      "sono": "sono",
424      "disney": "disney",
425      "high chair": "high chair",
426      "tv": "tv",
427      "hdtv": "tv", # Note: TV combines with HDTV
428      "hbo": "hbo",
429      "netflix": "netflix",
430      "pool": "pool",
431      "conditioner": "conditioner",
432      "soap": "soap",
433      "iron": "iron",
```

```
434        "sony": "sony",
435        "sound system": "sound system",
436        "heating": " radiant heating",
437        "stainless": "stainless",
438        "washer": "washer",
439        "dryer": "dryer",
440        "free parking": "free parking",
441        "baby monitor": "baby monitor",
442        "baby bath": "baby bath",
443        "body wash": "body wash",
444        "changing table": "changing table",
445        "books and toys": "books and toys",
446        "clothing storage": "clothing storage",
447        "exercise equipment": "exercise equipment",
448        "carport": "carport",
449        "free residential garage": "free residential garage",
450        "game console": "game console",
451        "fireplace": "fireplace",
452        "paid parking garage": "paid parking garage",
453        "paid parking lot off": "paid parking lot off",
454        "paid parking lot on": "paid parking lot on",
455        "paid parking on premises": "paid parking on premises",
456        "paid valet parking": "paid valet parking",
457        "hot water": "hot water",
458        "free driveway": "free driveway",
459
460 }
461
462 # List to store new aggregated amenity counts
463 aggregated_amenity_counts = {}
464
465 # Loop through amenities and aggregate based on keywords
466 for amenity, count in amenity_counts:
467     found = False
468     for keyword, new_amenity in combine_keywords.items():
469         if keyword in amenity.lower():
470             # Special handling for "oven" to ensure it doesn't get overshadowed by "toaster"
471             if keyword == "oven" and "toaster" in amenity.lower():
472                 continue
473             # Special handling for "hot water" to ensure "hot water kettle" doesn't get combined
474             if keyword == "hot water" and "kettle" in amenity.lower():
475                 continue
476             if new_amenity not in aggregated_amenity_counts:
477                 aggregated_amenity_counts[new_amenity] = 0
478             aggregated_amenity_counts[new_amenity] += count
479             found = True
480             break
481     if not found:
482         if amenity not in aggregated_amenity_counts:
483             aggregated_amenity_counts[amenity] = 0
484         aggregated_amenity_counts[amenity] += count
485
486 # Convert dictionary to list of tuples and sort
487 sorted_amenity_counts = sorted(aggregated_amenity_counts.items(), key=lambda x: x[1], reverse=True)
488
489 # Save sorted_amenity_counts to csv
```

```
490 df = pd.DataFrame(sorted_amenity_counts, columns=["Amenity", "Count"])
491 df.to_csv('new_amenity_counts.csv', index=False)
492
493 # Assuming nashvilleDF has a column 'amenities' that is either a string or list
494
495 # Ensure that the 'amenities' column is a list
496 nashvilleDF['amenities'] = nashvilleDF['amenities'].apply(lambda x: x if isinstance(x, list) else x.split(','))
497
498 # Define a function to aggregate amenities based on combine_keywords
499 def aggregate_amenities(amenities_list):
500     aggregated_list = []
501     for amenity in amenities_list:
502         found = False
503         for keyword, new_amenity in combine_keywords.items():
504             if keyword in amenity.lower():
505                 aggregated_list.append(new_amenity)
506                 found = True
507                 break
508         if not found:
509             aggregated_list.append(amenity)
510     return aggregated_list
511
512 # Apply the aggregate_amenities function to the amenities column
513 nashvilleDF['aggregated_amenities'] = nashvilleDF['amenities'].apply(aggregate_amenities)
514
515 # Create dummy variables
516 amenities_dummies = pd.get_dummies(nashvilleDF['aggregated_amenities'].apply(pd.Series).stack()).sum(level=0)
517
518 # Join the dummy variables to the original DataFrame and drop the amenities columns
519 nashvilleDF = nashvilleDF.join(amenities_dummies)
520 nashvilleDF = nashvilleDF.drop(["amenities", "aggregated_amenities"], axis=1)
521
522
523
524
525 # Convert columns to numeric and handle NaN values
526 columns_to_convert = [
527     'review_scores_rating', 'review_scores_accuracy', 'review_scores_cleanliness',
528     'review_scores_checkin', 'review_scores_communication', 'review_scores_location',
529     'review_scores_value'
530 ]
531
532 for col in columns_to_convert:
533     nashvilleDF[col] = pd.to_numeric(nashvilleDF[col], errors='coerce')
534     nashvilleDF[col].fillna(0, inplace=True)
535
536
537 # Define the size of the entire figure. (width, height)
538 fig, (ax1, ax2, ax3, ax4, ax5, ax6, ax7) = plt.subplots(7, 1, figsize=(10, 14))
539
540 ax1.hist(nashvilleDF.review_scores_rating, color='blue')
541 ax1.set_title('Rating')
542 ax2.hist(nashvilleDF.review_scores_accuracy, color='blue')
543 ax2.set_title('Accuracy')
544 ax3.hist(nashvilleDF.review_scores_cleanliness, color='blue')
545 ax3.set_title('Cleanliness')
```

```
546  ax4.hist(nashvilleDF.review_scores_checkin, color='blue')
547  ax4.set_title('Check In')
548  ax5.hist(nashvilleDF.review_scores_communication, color='blue')
549  ax5.set_title('Communication')
550  ax6.hist(nashvilleDF.review_scores_location, color='blue')
551  ax6.set_title('Location')
552  ax7.hist(nashvilleDF.review_scores_value, color='blue')
553  ax7.set_title('Value')
554
555  # Adjust the spacing between subplots
556  plt.subplots_adjust(hspace=0.5)
557
558  plt.tight_layout()  # This ensures that the titles and plots don't overlap
559  plt.show()
560
561
562  def findOutliers(df, column, lim_scalar=4):
563      """
564      Returns outliers above the max limit for a column in a dataframe
565      Adjust outlier cutoff to q75 + 4*iqr to include more data
566      ---
567      input: DataFrame, column(series),lim_scalar(float)
568      output: DataFrame
569      """
570      q25, q50, q75 = df[column].quantile(q=[0.25, 0.5, 0.75])
571      iqr = q75 - q25
572      # max limits to be considered an outlier
573      max_ = q75 + lim_scalar * iqr
574      # identify the points
575      outlier_mask = [True if x > max_ else False for x in df[column]]
576      print('{} outliers found out of {} data points, {}% of the data. {} is the max'.format(
577          sum(outlier_mask), len(df[column]),
578          100 * (sum(outlier_mask) / len(df[column])),max_))
579      return outlier_mask
580  #average number of bathrooms in nashvilleDF
581  # Bathrooms
582  # Fill na with 1
583  nashvilleDF.bathrooms.fillna(1,inplace=True)
584  nashvilleDF.bathrooms.mean()
585  # Look at the distribution of the bathrooms column
586  plt.figure(figsize=(3,3))
587  plt.boxplot(nashvilleDF.bathrooms)
588  sns.despine()
589  plt.title('bathrooms distribution');
590
591
592
593  # Remove bathroom outliers
594  nashvilleDF = nashvilleDF[np.logical_not(
595      findOutliers(nashvilleDF, 'bathrooms'))]
596  # Look at the distribution of the bathrooms column
597  plt.figure(figsize=(3,3))
598  plt.boxplot(nashvilleDF.bathrooms)
599  sns.despine()
600  plt.title('bathrooms distribution')
601
```

```
602
603
604  # Convert non-numeric values in the 'bedrooms' column to NaN
605  nashvilleDF['bedrooms'] = pd.to_numeric(nashvilleDF['bedrooms'], errors='coerce')
606
607  # Replace NaN values with 1
608  nashvilleDF['bedrooms'].fillna(1, inplace=True)
609
610  # Plot the boxplot for the cleaned 'bedrooms' column
611  plt.figure(figsize=(3,3))
612  plt.boxplot(nashvilleDF['bedrooms'])
613  sns.despine()
614  plt.title('beds distribution')
615  plt.show()
616
617
618
619
620  # Remove bedroom outliers
621  nashvilleDF = nashvilleDF[np.logical_not(
622      findOutliers(nashvilleDF, 'bedrooms',lim_scalar=6))]
623
624
625  # Plot the boxplot for the cleaned 'bedrooms' column
626  plt.figure(figsize=(3,3))
627  plt.boxplot(nashvilleDF['bedrooms'])
628  sns.despine()
629  plt.title('beds distribution')
630  plt.show()
631
632
633  # Plot distribution of the number of bathrooms
634  plt.figure(figsize=(10, 6))
635  plt.hist(nashvilleDF['bathrooms'], bins=20, color='blue', edgecolor='black')
636  plt.title('Distribution of Number of Bathrooms')
637  plt.xlabel('Number of Bathrooms')
638  plt.ylabel('Number of Listings')
639  plt.grid(axis='y')
640  plt.show()
641
642  # Plot distribution of the number of bathrooms
643  plt.figure(figsize=(10, 6))
644  plt.hist(nashvilleDF['bedrooms'], bins=20, color='blue', edgecolor='black')
645  plt.title('Distribution of Number of Bathrooms')
646  plt.xlabel('Number of Bedrooms')
647  plt.ylabel('Number of Listings')
648  plt.grid(axis='y')
649  plt.show()
650
651
652
653  # Calculate estimated number of bookings
654  nashvilleDF['est_bookings'] = nashvilleDF.number_of_reviews * 2
655
656  # Replace entries where unit is brand new with est_bookings = 1
657  nashvilleDF['est_bookings'] = [
```

```
658         1 if nashvilleDF.first_review.iloc[idx] == nashvilleDF.last_review.iloc[idx] else x
659         for idx, x in enumerate(nashvilleDF.est_bookings)
660     ]
661     # Convert 'minimum_nights' to numeric and handle NaN values (assuming a default value of 1)
662     nashvilleDF['minimum_nights'] = pd.to_numeric(nashvilleDF['minimum_nights'], errors='coerce')
663     nashvilleDF['minimum_nights'].fillna(1, inplace=True)
664
665     # Calculate estimated number of nights booked per year
666     # Use 3 days as the average length of a stay
667     # Unless the minimum number of days is greater than 3, then use that number
668     nashvilleDF['est_booked_nights_per_year'] = [
669         3 if x < 3 else x
670         for x in nashvilleDF.minimum_nights  # avg stay length
671     ] * nashvilleDF.reviews_per_month * 2 * 12
672     # Calculate estimated number of nights booked
673     # Use 3 days as the average length of a stay
674     # Unless the minimum number of days is greater than 3, then use that number
675
676     nashvilleDF['est_booked_nights'] = (
677         [
678             3 if x < 3 else x
679             for x in nashvilleDF.minimum_nights
680         ] *  # avg stay length
681         nashvilleDF['est_bookings'])
682     # Convert 'est_booked_nights' to numeric and handle NaN values
683     nashvilleDF['est_booked_nights'] = pd.to_numeric(nashvilleDF['est_booked_nights'], errors='coerce')
684     nashvilleDF['est_booked_nights'].fillna(0, inplace=True)
685
686     # Occupancy Rate = total_booked_nights / total_available_nights
687     nashvilleDF['occupancy_rate'] = nashvilleDF['est_booked_nights'] / (
688         (nashvilleDF.last_review - nashvilleDF.first_review).dt.days + 1)
689
690     # The next line seems redundant and is the same as the previous one. Consider removing it.
691     # Occupancy Rate = total_booked_nights / total_available_nights
692     nashvilleDF['occupancy_rate'] = nashvilleDF['est_booked_nights'] / (
693         (nashvilleDF.last_review - nashvilleDF.first_review).dt.days + 1)
694     # Convert 'availability_365' to numeric and handle NaN values
695     nashvilleDF['availability_365'] = pd.to_numeric(nashvilleDF['availability_365'], errors='coerce')
696     nashvilleDF['availability_365'].fillna(0, inplace=True)
697
698
699     # Convert 'est_booked_nights_per_year' to numeric and handle NaN values
700     nashvilleDF['est_booked_nights_per_year'] = pd.to_numeric(nashvilleDF['est_booked_nights_per_year'], errors='coerce')
701     nashvilleDF['est_booked_nights_per_year'].fillna(0, inplace=True)
702
703     # Calculate occupancy rate
704     #nashvilleDF['occupancy_rate2'] = nashvilleDF['est_booked_nights_per_year'] / (nashvilleDF['availability_365'] + 1)
705     Getting rid of data that isnt numerical or helpful.
706
707     columns_to_drop = [
708         "listing_url",
709         "scrape_id",
710         "source",
711         "name",
712         "description",
713         "neighborhood_overview",
```

```
714      "picture_url",
715      "host_id",
716      "host_url",
717      "host_name",
718      "host_location",
719      "host_about",
720      "host_thumbnail_url",
721      "host_picture_url",
722      "host_neighbourhood",
723      "host_verifications",
724      "neighbourhood_cleansed",
725      "neighbourhood_group_cleansed",
726      "latitude",
727      "longitude",
728      'property_type',
729      'room_type',
730      'est_booked_nights_per_year',
731      'est_booked_nights',
732      'occupancy_rate',
733      'host_listings_count',
734      'est_bookings',
735      'calculated_host_listings_count',
736      'calculated_host_listings_count_private_rooms',
737      'minimum_nights_avg_ntm',
738      'maximum_nights_avg_ntm',
739      'calculated_host_listings_count_entire_homes',
740      'host_total_listings_count',
741      'availability_90',
742      'availability_60', 'minimum_nights',
743      'maximum_minimum_nights', 'review_scores_communication',
744      'smaller', 'ge', 'irish spring', 'lotion','suave', 'dove', 'dove anti-stress moisturizing cream bar','dr teals',
             'smaller fridge',
745      'organic', 'olympic-sized','number_of_reviews_l30d', 'number_of_reviews_ltm',
746      '5-10 years old', '2-5 years old',
747
748
749 ]
750 nashvilleDF['price'] = nashvilleDF['price'].str.replace('$', '').str.replace(',', '').astype(float)
751
752 nashvilleDF = nashvilleDF.drop(columns=columns_to_drop, errors='ignore')
753
754
755 # Convert the 'price' column to numeric (assuming it contains strings representing numbers)
756 nashvilleDF['price'] = pd.to_numeric(nashvilleDF['price'], errors='coerce')
757
758 try:
759      # Calculate the required statistics
760      total_entries = len(nashvilleDF['price'])
761      entries_500_or_less = len(nashvilleDF[nashvilleDF['price'] <= 500])
762      entries_over_500 = len(nashvilleDF[nashvilleDF['price'] > 1000])
763
764      results = total_entries, entries_500_or_less, entries_over_500
765 except Exception as e:
766      error_message = str(e)
767      results = None
768
```

```
769  results, error_message if results is None else None

770

771

772  #drop first_review  last_review last_scraped  host_since
773  nashvilleDF.drop(['first_review','last_review','last_scraped','host_since',
774                    'calendar_updated','calendar_last_scraped','license'],axis=1,inplace=True)

775

776  #host_acceptance_rate has % on it, drop the % from it
777  nashvilleDF['host_acceptance_rate'] = nashvilleDF['host_acceptance_rate'].str.replace('%', '').astype(float)
778  missing_data = nashvilleDF.isnull().sum()
779  print(missing_data[missing_data > 0])

780

781

782

783  nashvilleDF = nashvilleDF.apply(pd.to_numeric, errors='coerce')

784

785  nashvilleDF["host_response_rate"] = nashvilleDF["host_response_rate"].fillna(nashvilleDF["host_response_rate"].median())
786  nashvilleDF["host_acceptance_rate"] =
787        nashvilleDF["host_acceptance_rate"].fillna(nashvilleDF["host_acceptance_rate"].median())

787

788  # Impute missing values for beds with the median
789  nashvilleDF["beds"] = nashvilleDF["beds"].fillna(nashvilleDF["beds"].median())

790

791  # Convert reviews_per_month to an integer column
792  nashvilleDF["reviews_per_month"] = nashvilleDF["reviews_per_month"].astype("float64")

793

794  # Impute missing values for reviews_per_month and occupancy_rate with the mean
795  nashvilleDF["reviews_per_month"] = nashvilleDF["reviews_per_month"].fillna(nashvilleDF["reviews_per_month"].mean())
796  #limit price between 200-300
797  nashvilleDF = nashvilleDF[nashvilleDF['price'] >= 000]

798

799

800  nashvilleDF = nashvilleDF[nashvilleDF['price'] <= 1000]
801  missing_data = nashvilleDF.isnull().sum()
802  print(missing_data[missing_data > 0])
803  Series([], dtype: int64)
804  #write nashvilleDF to a csv file
805  nashvilleDF.to_csv('nashvilleDF.csv', index=False)
806  plt.figure(figsize=(10, 6))
807  sns.histplot(nashvilleDF['price'], bins=50, kde=True)
808  plt.title('Distribution of Price')
809  plt.xlabel('Price')
810  plt.ylabel('Frequency')
811  plt.show()

812

813

814  # Split data into features (X) and target (y)
815  X = nashvilleDF.drop('price', axis=1)
816  y = nashvilleDF['price']

817

818  # Split data into training and testing sets
819  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

820

821  # Initialize linear regression model
822  model = LinearRegression()

823
```

```
824  # Train the model
825  model.fit(X_train, y_train)
826
827  # Predict on test set
828  y_pred = model.predict(X_test)
829  # Calculate and print R^2 score
830  r2 = r2_score(y_test, y_pred)
831  print(f"R^2 Score: {r2}")
832
833
834  import pandas as pd
835  import numpy as np
836  import matplotlib.pyplot as plt
837  from sklearn.linear_model import LinearRegression, SGDRegressor
838  from sklearn.tree import DecisionTreeRegressor
839  from sklearn.ensemble import GradientBoostingRegressor, RandomForestRegressor, HistGradientBoostingRegressor
840  from sklearn.model_selection import train_test_split
841  from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
842  import xgboost as xgb
843
844  subsetdf = nashvilleDF.copy()
845  df = subsetdf.copy()
846
847  def prepare_data(df, target_col='price'):
848      X = df.drop(target_col, axis=1)
849      y = df[target_col]
850      return train_test_split(X, y, test_size=0.2, random_state=42)
851
852  def create_models():
853      return {
854          'GradientBoosting': GradientBoostingRegressor(),
855          'RandomForest': RandomForestRegressor(),
856          'Linear': LinearRegression(),
857          'HistGradientBoosting': HistGradientBoostingRegressor(),
858          'DecisionTree': DecisionTreeRegressor(),
859          'XGBoost': xgb.XGBRegressor(),
860      }
861
862  def fit_models(X_train, y_train, models):
863      for name, model in models.items():
864          model.fit(X_train.astype(float), y_train.astype(float))
865      return models
866
867  def evaluate_models(X_test, y_test, models):
868      results = []
869      for name, model in models.items():
870          y_pred = model.predict(X_test.astype(float))
871
872          r2 = r2_score(y_test.astype(float), y_pred)
873          mse = mean_squared_error(y_test.astype(float), y_pred)
874          rmse = np.sqrt(mse)
875          mae = mean_absolute_error(y_test.astype(float), y_pred)
876
877          results.append({
878              'Model': name,
879              'R2': r2,
```

```
880              'RMSE': rmse,
881              'MSE': mse,
882              'MAE': mae
883          })
884      return pd.DataFrame(results)
885
886 def plot_model_performance(results_df):
887      fig, axes = plt.subplots(2, 2, figsize=(10, 8))
888      fig.suptitle('Model Performance Comparison')
889
890      metrics = ['R2', 'RMSE', 'MSE', 'MAE']
891
892      for idx, metric in enumerate(metrics):
893          ax = axes[idx // 2, idx % 2]
894          results_df.plot(x='Model', y=metric, kind='bar', ax=ax, legend=None)
895          ax.set_ylabel(metric)
896          ax.set_ylim(bottom=0)
897          ax.set_xticklabels(results_df['Model'], rotation=45, ha='right')
898
899      plt.tight_layout(rect=[0, 0, 1, 0.96])
900      plt.show()
901
902 def main(df):
903      X_train, X_test, y_train, y_test = prepare_data(df)
904      models = create_models()
905      models = fit_models(X_train, y_train, models)
906      results_df = evaluate_models(X_test, y_test, models)
907      plot_model_performance(results_df)
908      table = pd.pivot_table(results_df, index='Model', values=['R2', 'RMSE', 'MSE', 'MAE'])
909      print(table)
910
911 main(df)
912
913
914
915
916
917
918
919
920 from sklearn.ensemble import HistGradientBoostingRegressor
921 from sklearn.inspection import permutation_importance
922 from sklearn.model_selection import KFold
923
924 # Initialize the HistGradientBoostingRegressor model
925 model = HistGradientBoostingRegressor()
926
927 # Specify the number of folds for k-fold cross-validation
928 num_folds = 10
929
930 # Initialize the KFold object
931 kf = KFold(n_splits=num_folds, shuffle=True, random_state=42)
932
933
934 # Initialize an array to store the accuracy scores
935 scores = np.zeros(num_folds)
```

```
936
937  # Initialize an array to store feature importances for each fold
938  # The number of columns in X is assumed to be X.shape[1]
939  feature_importances = np.zeros((num_folds, X.shape[1]))
940
941  # Loop through each fold
942  for fold, (train_index, test_index) in enumerate(kf.split(X)):
943      # Split the data into training and test sets
944      X_train, y_train = X.iloc[train_index], y.iloc[train_index]
945      X_test, y_test = X.iloc[test_index], y.iloc[test_index]
946
947      # Train your model on the training set
948      model.fit(X_train, y_train)
949
950      # Evaluate your model on the test set
951      scores[fold] = model.score(X_test, y_test)
952
953      # Get feature importances using permutation importance
954      perm_importance = permutation_importance(model, X_test, y_test, n_repeats=10, random_state=42)
955      feature_importances[fold] = perm_importance.importances_mean
956
957      # Print the results
958      print(f'Fold {fold+1}: Accuracy = {scores[fold]}')
959
960  # Compute and print the average accuracy across all folds
961  print(f'Average Accuracy: {np.mean(scores)}')
962
963  # Compute the average feature importance across all folds
964  avg_feature_importances = feature_importances.mean(axis=0)
965
966  # Get the top 20 features based on average importance
967  top_features = np.argsort(avg_feature_importances)[-20:]
968  top_feature_names = X.columns[top_features]
969  top_feature_importances = avg_feature_importances[top_features]
970
971  top_feature_names, top_feature_importances
972
973
974
975
976
977
978  import pandas as pd
979  import numpy as np
980  from sklearn.experimental import enable_hist_gradient_boosting
981  from sklearn.ensemble import HistGradientBoostingRegressor
982  from sklearn.inspection import permutation_importance
983  from collections import defaultdict
984
985  # Define function to get feature importances
986  def get_feature_importance(data, features):
987      X = data[features]
988      y = data['price']
989
990      model = HistGradientBoostingRegressor()
991      model.fit(X, y)
```

```
992
993      # Using permutation importance to get feature importances
994      result = permutation_importance(model, X, y, n_repeats=10, random_state=42)
995      importances = result.importances_mean
996
997      return dict(zip(features, importances))
998
999  # Iteratively get feature importances for random feature subsets
1000 iterations = 100  # Number of iterations
1001 num_features = 20
1002
1003 feature_columns = [col for col in nashvilleDF.columns if col != 'price']
1004
1005 collective_importances = defaultdict(float)
1006
1007 for _ in range(iterations):
1008     selected_features = np.random.choice(feature_columns, num_features, replace=False)
1009
1010     importances = get_feature_importance(nashvilleDF, selected_features)
1011
1012     for feature, importance in importances.items():
1013         collective_importances[feature] += importance
1014
1015 # Sort features by their collective importance
1016 sorted_features = sorted(collective_importances.keys(), key=lambda x: collective_importances[x], reverse=True)
1017
1018 print(sorted_features)
1019
1020 from sklearn.experimental import enable_hist_gradient_boosting
1021 from sklearn.ensemble import HistGradientBoostingRegressor
1022 from sklearn.inspection import permutation_importance
1023
1024 # Split the data into features and target
1025 X = nashvilleDF.drop('price', axis=1)
1026 y = nashvilleDF['price']
1027
1028 # Initialize the model
1029 model = HistGradientBoostingRegressor()
1030
1031 # Fit the model
1032 model.fit(X, y)
1033
1034 # Get feature importances using permutation importance
1035 perm_importance = permutation_importance(model, X, y, n_repeats=30, random_state=42)
1036
1037 # Get the top 20 features
1038 top_20_features = sorted(zip(perm_importance.importances_mean, X.columns), reverse=True)[:20]
1039 top_20_features
1040
1041
1042
1043 # Categorizing amenities into defined buckets
1044 kitchen_amenities = [
1045     "coffee", "kitchen", "refrigerator", "dishes and silverware", "microwave",
1046     "oven", "stove", "cooking basics", "freezer", "wine glasses", "toaster",
1047     "hot water kettle", "grill", "blender", "barbecue utensils", "rice maker",
```

```
1048        "nespresso", "espresso machine", "french press", "bread maker", "dining table",
1049        "baking sheet", "gas", "we have a countertop burner"
1050 ]
1051
1052 cleaning_amenities = [
1053        "dryer", "washer", "shampoo", "iron", "hot water", "soap", "cleaning products",
1054        "shower gel", "conditioner", "cleaning available during stay", "body wash",
1055        "lotion", "suave", "irish spring", "la botegga", "dove anti-stress moisturizing cream bar"
1056 ]
1057
1058 safety_amenities = [
1059        "smoke alarm", "carbon monoxide alarm", "fire extinguisher", "first aid kit",
1060        "keypad", "smart lock", "lockbox", "camera", "lock on bedroom door", "safe",
1061        "baby safety gates", "outlet covers", "table corner guards", "window guards",
1062        "ev charger", "ev charger - level 2", "ev charger - level 1", "tesla only"
1063 ]
1064
1065 household_amenities = [
1066        "radiant heating", "essentials", "hangers", "self check-in", "bed linens",
1067        "free parking", "air conditioning", "private entrance", "extra pillows and blankets",
1068        "dedicated workspace", "free street parking", "long term stays allowed", "clothing storage",
1069        "ceiling fan", "central air conditioning", "private patio or balcony", "room-darkening shades",
1070        "luggage dropoff allowed", "single level home", "elevator", "fireplace",
1071        "books and reading material", "laundromat nearby", "board games", "backyard",
1072        "outdoor furniture", "patio or balcony", "outdoor dining area", "gym",
1073        "fire pit", "high chair", "crib", "pool", "hot tub", "breakfast", "resort access",
1074        "city skyline view", "exercise equipment", "portable fans", "childrenu2019s dinnerware",
1075        "building staff", "babysitter recommendations", "free driveway", "closet",
1076        "shared patio or balcony", "record player", "trash compactor", "sun loungers",
1077        "closet", "wardrobe", "portable heater", "bidet", "changing table", "bikes",
1078        "hammock", "2-5 years old", "5-10 years old", "and 10 years old", "and 5-10 years old",
1079        "rooftop", "saltwater", "infinity", "beach essentials", "mosquito net", "heated",
1080        "dvd player", "ev charger - level 1", "private sauna", "shared sauna", "window ac unit",
1081        "host greets you", "piano", "baby monitor", "ping pong table", "private living room",
1082        "portable air conditioning", "smoking allowed", "baby bath", "beach access", "ski-in/ski-out",
1083        "wood-burning", "outdoor shower", "ac - split type ductless system", "window ac unit",
1084        "all natural", "stainless", "portable", "sono", "but good size not standard large size",
1085        "smaller", "some type", "two racks", "organic", "olympic-sized", "smaller fridge"
1086 ]
1087
1088 bedroom_amenities = [
1089        "hangers", "bed linens", "extra pillows and blankets", "room-darkening shades",
1090        "crib", "single level home", "closet", "wardrobe", "and dresser", "changing table"
1091 ]
1092
1093 electronics_amenities = [
1094        "wifi", "tv", "ethernet connection", "cable", "roku", "hulu", "disney", "hbo",
1095        "netflix", "amazon prime video", "bluetooth", "game console", "xbox 360", "and xbox one",
1096        "sound system", "vizio soundbar", "chromecast"
1097 ]
1098
1099 extra_spaces_amenities = [
1100        "patio or balcony", "pool", "garden or backyard", "hot tub", "BBQ grill",
1101        "private living room", "gym", "fire pit", "outdoor dining area", "backyard",
1102        "outdoor furniture", "private patio or balcony", "gym", "fireplace", "rooftop",
1103        "beach access", "ski-in/ski-out", "outdoor shower", "resort access", "waterfront",
```

```
1104         "resort view", "harbor view", "beach view", "sea view", "ocean view",
1105         "lake access", "mountain view", "valley view", "golf course view", "lake view",
1106         "marina view", "public or shared beach access", "canal view", "boat slip", "bay view",
1107         "park view"
1108  ]
1109
1110  review_amenities = [
1111         "number_of_reviews", "number_of_reviews_ltm", "number_of_reviews_l30d", "review_scores_rating",
1112         "review_scores_accuracy", "review_scores_cleanliness", "review_scores_checkin", "review_scores_communication",
1113         "review_scores_location", "review_scores_value"
1114  ]
1115
1116
1117  # Return the categorized amenities
1118  categorized_amenities = {
1119         "Kitchen Amenities": kitchen_amenities,
1120         "Cleaning Amenities": cleaning_amenities,
1121         "Safety Amenities": safety_amenities,
1122         "Household Amenities": household_amenities,
1123         "Bedroom Amenities": bedroom_amenities,
1124         "Electronics Amenities": electronics_amenities,
1125         "Extra Spaces Amenities": extra_spaces_amenities,
1126         "Review Amenities": review_amenities
1127
1128  }
1129  import networkx as nx
1130  import matplotlib.pyplot as plt
1131
1132  # Create a new directed graph
1133  G = nx.DiGraph()
1134
1135  # Add nodes for each amenity group
1136  amenity_groups_names = [
1137         "Kitchen Amenities", "Cleaning Amenities", "Safety Amenities",
1138         "Household Amenities", "Bedroom Amenities", "Electronics Amenities",
1139         "Extra Spaces Amenities", "Review Amenities"
1140  ]
1141  for group in amenity_groups_names:
1142         G.add_node(group)
1143
1144  # Add edges from each amenity group to 'price'
1145  for group in amenity_groups_names:
1146         G.add_edge(group, 'price')
1147
1148  # Plot the DAG
1149  plt.figure(figsize=(10, 6))
1150
1151  # Use shell_layout instead of spring_layout
1152  shell_nestings = [amenity_groups_names, ['price']]  # Define the concentric circles
1153  pos = nx.shell_layout(G, nlist=shell_nestings)
1154
1155  nx.draw_networkx(G, pos, with_labels=True, node_size=4000, node_color='skyblue', font_size=10, font_weight='bold',
1156         width=2, edge_color='gray')
1156  plt.title("Causal DAG representing the effect of amenity groups on price")
1157  plt.show()
1158
```

```
1159
1160
1161  from sklearn.model_selection import KFold
1162  import numpy as np
1163
1164  # Load your data into a pandas DataFrame
1165  data = nashvilleDF.copy()
1166
1167  # Split your data into features and target
1168  X = data.drop('price', axis=1)
1169  y = data['price']
1170
1171  # Define the number of folds
1172  num_folds = 10
1173
1174  # Initialize the cross-validation method
1175  kf = KFold(n_splits=num_folds)
1176
1177  # Initialize a linear regression model
1178  model = HistGradientBoostingRegressor()
1179
1180  # Initialize an array to store the accuracy scores
1181  scores = np.zeros(num_folds)
1182
1183  # Loop through each fold
1184  for fold, (train_index, test_index) in enumerate(kf.split(X)):
1185
1186      # Split the data into training and test sets
1187      X_train, y_train = X.iloc[train_index], y.iloc[train_index]
1188      X_test, y_test = X.iloc[test_index], y.iloc[test_index]
1189
1190      # Train your model on the training set
1191      model.fit(X_train, y_train)
1192
1193      # Evaluate your model on the test set
1194      scores[fold] = model.score(X_test, y_test)
1195
1196      # Print the results
1197      print(f'Fold {fold+1}: Accuracy = {scores[fold]}')
1198
1199  # Compute and print the average accuracy across all folds
1200  print(f'Average Accuracy: {np.mean(scores)}')
1201
1202
1203  model1_plot = nashvilleDF.groupby(['accommodates', 'bedrooms']).size().reset_index().pivot(columns='bedrooms',
1204                                                                                    index='accommodates', values=0)
1205
1206  model1_plot.plot(kind='bar', stacked=True, figsize=(12,5))
1207  plt.title("Accomodates and Number of Bedrooms", size=20)
1208  plt.xlabel("Accomodates", size=15)
1209
1210
1211
1212
1213  # Filter rows where 'wifi' column equals 1
1214  wifi_1_rows = nashvilleDF[nashvilleDF['wifi'] == 1]
```

71

```
1215
1216 # Calculate the percentage
1217 percent_with_wifi_1 = (len(wifi_1_rows) / len(nashvilleDF)) * 100
1218
1219 print(f"Percentage of rows with 'wifi' = 1 in nashvilleDF: {percent_with_wifi_1:.2f}%")
1220
1221
1222 # Filter rows where 'wifi' column equals 1
1223 resort_1_rows = nashvilleDF[nashvilleDF['resort access'] == 1]
1224
1225 # Calculate the percentage
1226 percent_with_resort_1 = (len(resort_1_rows) / len(nashvilleDF)) * 100
1227
1228 print(f"Percentage of rows with 'wifi' = 1 in nashvilleDF: {percent_with_resort_1:.2f}%")
1229
1230 from sklearn.ensemble import HistGradientBoostingRegressor
1231 from sklearn.inspection import permutation_importance
1232 from sklearn.model_selection import train_test_split, KFold
1233 from joblib import Parallel, delayed
1234 import numpy as np
1235
1236 def process_fold(fold, train_index, test_index, X, y):
1237     model = HistGradientBoostingRegressor()
1238     X_train, y_train = X.iloc[train_index], y.iloc[train_index]
1239     X_test, y_test = X.iloc[test_index], y.iloc[test_index]
1240     model.fit(X_train, y_train)
1241     score = model.score(X_test, y_test)
1242     perm_importance = permutation_importance(model, X_test, y_test, n_repeats=10, random_state=42)
1243     sorted_idx = perm_importance.importances_mean.argsort()[::-1]
1244     sorted_features = [(X.columns[i], perm_importance.importances_mean[i]) for i in sorted_idx]
1245     return score, perm_importance.importances_mean, sorted_features
1246
1247 # Assuming nashvilleDF is defined
1248 X = nashvilleDF.drop(['price', 'id'], axis=1)
1249 y = nashvilleDF['price']
1250 X.columns = X.columns.astype(str)
1251
1252 num_folds = 10
1253 kf = KFold(n_splits=num_folds, shuffle=True, random_state=42)
1254 scores = np.zeros(num_folds)
1255 feature_importances = np.zeros((num_folds, X.shape[1]))
1256 results = Parallel(n_jobs=-1)(delayed(process_fold)(fold, train_index, test_index, X, y)
1257                              for fold, (train_index, test_index) in enumerate(kf.split(X)))
1258
1259 accuracy_scores = []
1260 for fold, (score, importance, sorted_features) in enumerate(results):
1261     scores[fold] = score
1262     feature_importances[fold] = importance
1263
1264     # Append the score to accuracy_scores list
1265     accuracy_scores.append(score)
1266
1267     print(f'Fold {fold+1}: Accuracy = {scores[fold]}')
1268     print("Features ranked by importance for fold:", fold+1)
1269     for name, imp in sorted_features:
1270         print(f"{name}: {imp}")
```

72

```
1271
1272  # Compute and print the average accuracy and standard deviation across all folds
1273  mean_accuracy = np.mean(accuracy_scores)
1274  std_accuracy = np.std(accuracy_scores)
1275  print(f'Average Accuracy: {mean_accuracy}')
1276  print(f'Standard Deviation of Accuracy: {std_accuracy}')
1277
1278  # Compute the average feature importance across all folds
1279  avg_feature_importances = feature_importances.mean(axis=0)
1280
1281  # Get the top 20 features based on average importance
1282  top_features = np.argsort(avg_feature_importances)[-20:]
1283  top_feature_names = X.columns[top_features]
1284  top_feature_importances = avg_feature_importances[top_features]
1285
1286  print("Top 20 Features and their Importances:")
1287  print(list(zip(top_feature_names, top_feature_importances)))
1288
1289
1290  all_top_features = []
1291  for _, _, sorted_features in results:
1292      top_20 = [feature[0] for feature in sorted_features[:20]]
1293      all_top_features.extend(top_20)
1294
1295  # Calculate the average rank for each feature
1296  feature_ranks = {feature: 0 for feature in set(all_top_features)}
1297  for feature in all_top_features:
1298      feature_ranks[feature] += (all_top_features.count(feature) - all_top_features.index(feature) - 1)
1299
1300  average_ranks = {feature: rank / all_top_features.count(feature) for feature, rank in feature_ranks.items()}
1301
1302  # Sort by average rank
1303  sorted_avg_ranks = sorted(average_ranks.items(), key=lambda x: x[1], reverse=True)
1304
1305  # Displaying the sorted average ranks
1306  for feature, rank in sorted_avg_ranks:
1307      print(f"{feature}: {rank}")
1308
1309  # Plotting the graph
1310  features = [item[0] for item in sorted_avg_ranks]
1311  ranks = [item[1] for item in sorted_avg_ranks]
1312
1313  plt.figure(figsize=(12, 10))
1314  plt.barh(features, ranks, color='skyblue')
1315  plt.xlabel('Average Rank')
1316  plt.ylabel('Feature Name')
1317  plt.title('Average Rank of Top Features')
1318  plt.gca().invert_yaxis()
1319  plt.tight_layout()
1320  plt.show()
```

## CAUSAL DISCOVERY

# B.1     Large Language Model Causal Discovery

```
1  !pip install langchain.agents
2  import os
3  from itertools import combinations
4
5  import numpy as np
6  from scipy import linalg
7  from scipy import stats
8
9  import matplotlib.pyplot as plt
10
11 from langchain.agents import load_tools, initialize_agent
12 from langchain.agents import AgentType
13
14 from langchain.chat_models import ChatOpenAI
15
16 from castle.common import GraphDAG
17 from castle.metrics import MetricsDAG
18 from castle.algorithms import PC
19
20 from castle.common.priori_knowledge import PrioriKnowledge
21
22
23 COLORS = [
24     '#00B0F0',
25     '#FF0000',
26     '#B0F000'
27 ]
28
29
30
31 def check_if_dag(A):
32     return np.trace(linalg.expm(A * A)) - A.shape[0] == 0
33
34
35
36
37 #Put in your OPENAI KEY HERE
38 os.environ['OPENAI_API_KEY'] = "sk-KEYASGFASGUASOGUHA"
39
40 all_vars = {
41     'altitude': 0,
42     'oxygen_density': 1,
```

```
43      'temperature': 2,
44      'risk_of_death': 3,
45      'mehendretex': 4
46  }
47
48
49  SAMPLE_SIZE = 1000
50
51  altitude = stats.halfnorm.rvs(scale=2000, size=SAMPLE_SIZE)
52  temperature = 25 - altitude / 100 + stats.norm.rvs(
53      loc=0,
54      scale=2,
55      size=SAMPLE_SIZE
56  )
57
58  mehendretex = stats.halfnorm.rvs(size=SAMPLE_SIZE)
59
60  oxygen_density = np.clip(
61      1 - altitude / 8000
62      - temperature / 50
63      + stats.norm.rvs(size=SAMPLE_SIZE) / 20,
64      0,
65      1)
66
67  risk_of_death = np.clip(
68      altitude / 20000
69      + np.abs(temperature) / 100
70      - oxygen_density / 5
71      - mehendretex / 5
72      + stats.norm.rvs(size=SAMPLE_SIZE) / 10,
73      0,
74      1
75  )
76
77
78
79  dataset = np.stack(
80      [
81          altitude,
82          oxygen_density,
83          temperature,
84          risk_of_death,
85          mehendretex
86      ]
87  ).T
88
89
90
91  true_dag = np.array(
92      [
93          [0, 1, 1, 1, 0],
94          [0, 0, 0, 1, 0],
95          [0, 1, 0, 1, 0],
96          [0, 0, 0, 0, 0],
97          [0, 0, 0, 1, 0]
98      ]
```

```
 99 )
100
101
102 # PC discovery without LLM assist
103 pc = PC(variant='stable')
104 pc.learn(dataset)
105
106 print(f'Is learned matrix a DAG: {check_if_dag(pc.causal_matrix)}')
107
108 # Vizualize
109 GraphDAG(
110     est_dag=pc.causal_matrix,
111     true_dag=true_dag)
112
113 plt.show()
114
115 # Compute metrics
116 metrics = MetricsDAG(
117     B_est=pc.causal_matrix,
118     B_true=true_dag)
119
120 print(metrics.metrics)
121
122 # Instantiate and encode priori knowledge
123 priori_knowledge = PrioriKnowledge(n_nodes=len(all_vars))
124
125 llm = ChatOpenAI(
126     temperature=0, # Temp == 0 => we want clear reasoning
127     model='gpt-4')#'gpt-3.5-turbo')
128
129
130 # Load tools
131 tools = load_tools(
132     [
133         "wikipedia"
134     ],
135     llm=llm)
136
137
138 # Instantiate the agent
139 agent = initialize_agent(
140     tools,
141     llm,
142     agent=AgentType.CHAT_ZERO_SHOT_REACT_DESCRIPTION,
143     handle_parsing_errors=True,
144     verbose=False)
145
146
147
148 def get_llm_info(llm, agent, var_1, var_2):
149
150     out = agent(f"Does {var_1} cause {var_2} or the other way around?\
151     We assume the following definition of causation:\
152     if we change A, B will also change.\
153     The relationship does not have to be linear or monotonic.\
154     We are interested in all types of causal relationships, including\
```

```python
155        partial and indirect relationships , given that our definition holds.\
156        ")
157
158        print(out)
159
160        pred = llm.predict(f'We assume the following definition of causation:\
161        if we change A, B will also change.\
162        Based on the following information: {out["output"]},\
163        print (0,1) if {var_1} causes {var_2},\
164        print (1, 0) if {var_2} causes {var_1}, print (0,0)\
165        if there is no causal relationship between {var_1} and {var_2}.\
166        Finally, print (-1, -1) if you don\'t know. Importantly, don\'t try to\
167        make up an answer if you don\'t know.')
168
169        print(pred)
170
171        return pred
172
173 for var_1, var_2 in combinations(all_vars.keys(), r=2):
174        print(var_1, var_2)
175        out = get_llm_info(llm, agent, var_1, var_2)
176        if out=='(0,1)':
177            priori_knowledge.add_required_edges(
178                [(all_vars[var_1], all_vars[var_2])]
179            )
180
181            priori_knowledge.add_forbidden_edges(
182                [(all_vars[var_2], all_vars[var_1])]
183            )
184
185        elif out=='(1,0)':
186            priori_knowledge.add_required_edges(
187                [(all_vars[var_2], all_vars[var_1])]
188            )
189            priori_knowledge.add_forbidden_edges(
190                [(all_vars[var_1], all_vars[var_2])]
191            )
192
193 print('\nLLM knowledge vs true DAG')
194 priori_dag = np.clip(priori_knowledge.matrix, 0, 1)
195
196 print(f'\nChecking if priori graph is a DAG: {check_if_dag(priori_dag)}')
197
198 GraphDAG(
199        est_dag=priori_dag,
200        true_dag=true_dag)
201
202 plt.show()
203
204 print('\nRunning PC')
205
206 # Instantiate the model with expert knowledge
207 pc_priori = PC(
208        priori_knowledge=priori_knowledge,
209        variant='stable'
210 )
```

```
211
212  # Learn
213  pc_priori.learn(dataset)
214
215  GraphDAG(
216      est_dag=pc_priori.causal_matrix,
217      true_dag=true_dag)
218
219  plt.show()
220
221  # Compute metrics
222  metrics = MetricsDAG(
223      B_est=pc_priori.causal_matrix,
224      B_true=true_dag)
225
226  print(metrics.metrics)
227
228  import numpy as np
229
230  # Define the order of variables
231  variables_order = ["accommodates", "price", "bedrooms", "beds", "bathrooms"]
232
233  # Initialize the adjacency matrix with zeros
234  n_vars = len(variables_order)
235  true_dag = np.zeros((n_vars, n_vars))
236
237  # Define the relationships as given in the digraph format
238  relationships = [
239      ("accommodates", "price"),
240      ("bedrooms", "accommodates"),
241      ("beds", "bedrooms"),
242      ("bathrooms", "bedrooms"),
243      ("bathrooms", "price")
244  ]
245
246  # Fill in the matrix based on the relationships
247  for source, target in relationships:
248      i = variables_order.index(source)
249      j = variables_order.index(target)
250      true_dag[i][j] = 1
251
252  print(true_dag)
253
254
255
256  import os
257  from itertools import combinations
258  import numpy as np
259  import pandas as pd
260  import matplotlib.pyplot as plt
261  from castle.common import GraphDAG
262  from castle.metrics import MetricsDAG
263  from castle.algorithms import PC
264
265  import os
266  from itertools import combinations
```

```python
267  import numpy as np
268  import pandas as pd
269  import matplotlib.pyplot as plt
270  from castle.common import GraphDAG
271  from castle.metrics import MetricsDAG
272  from castle.algorithms import PC
273  from castle.common.priori_knowledge import PrioriKnowledge
274
275  # Load the dataset
276  file_path = r"C:\Users\nstep\TSU\SeniorProject\df_selected1.csv"
277  nashvilleDF = pd.read_csv(file_path)
278  data_subset = nashvilleDF[['accommodates', 'price','bedrooms','beds','bathrooms']].dropna().values
279
280  # Define the order of variables and create the true_dag
281  variables_order = ["accommodates", "price", "bedrooms", "beds", "bathrooms"]
282  all_vars = {var: idx for idx, var in enumerate(variables_order)}
283  n_vars = len(variables_order)
284  true_dag = np.zeros((n_vars, n_vars))
285  relationships = [
286      ("accommodates", "price"),
287      ("bedrooms", "accommodates"),
288      ("beds", "bedrooms"),
289      ("bathrooms", "bedrooms"),
290      ("bathrooms", "price")
291  ]
292  for source, target in relationships:
293      i = variables_order.index(source)
294      j = variables_order.index(target)
295      true_dag[i][j] = 1
296
297  # Run PC discovery without LLM assist
298  pc = PC(variant='stable')
299  pc.learn(data_subset)
300
301  # Visualize the result
302  GraphDAG(
303      est_dag=pc.causal_matrix,
304      true_dag=true_dag)
305
306  plt.show()
307
308  # Compute metrics
309  metrics = MetricsDAG(
310      B_est=pc.causal_matrix,
311      B_true=true_dag)
312  print(metrics.metrics)
313
314  # Instantiate and encode priori knowledge
315  priori_knowledge = PrioriKnowledge(n_nodes=len(all_vars))
316
317  # Instantiate the GPT agent (You'll need to provide the appropriate setup and imports for the LLM)
318  # This step is based on the code you provided
319  llm = ChatOpenAI(
320      temperature=0, # Temp == 0 => we want clear reasoning
321      model='gpt-4')#'gpt-3.5-turbo')
322
```

```
323 tools = load_tools(
324     ["wikipedia"],
325     llm=llm)
326
327 agent = initialize_agent(
328     tools,
329     llm,
330     agent=AgentType.CHAT_ZERO_SHOT_REACT_DESCRIPTION,
331     handle_parsing_errors=True,
332     verbose=False)
333
334 def get_llm_info(llm, agent, var_1, var_2):
335
336     out = agent(f"Does {var_1} cause {var_2} or the other way around?\
337     We assume the following definition of causation:\
338     if we change A, B will also change.\
339     The relationship does not have to be linear or monotonic.\
340     We are interested in all types of causal relationships, including\
341     partial and indirect relationships, given that our definition holds.\
342     ")
343
344     print(out)
345
346     pred = llm.predict(f'We assume the following definition of causation:\
347     if we change A, B will also change.\
348     Based on the following information: {out["output"]},\
349     print (0,1) if {var_1} causes {var_2},\
350     print (1, 0) if {var_2} causes {var_1}, print (0,0)\
351     if there is no causal relationship between {var_1} and {var_2}.\
352     Finally, print (-1, -1) if you don\'t know. Importantly, don\'t try to\
353     make up an answer if you don\'t know.')
354
355     print(pred)
356
357     return pred
358
359 # Add priori knowledge from the LLM
360 for var_1, var_2 in combinations(all_vars.keys(), r=2):
361     print(var_1, var_2)
362     out = get_llm_info(llm, agent, var_1, var_2)
363     if out=='(0,1)':
364         priori_knowledge.add_required_edges([(all_vars[var_1], all_vars[var_2])])
365         priori_knowledge.add_forbidden_edges([(all_vars[var_2], all_vars[var_1])])
366     elif out=='(1,0)':
367         priori_knowledge.add_required_edges([(all_vars[var_2], all_vars[var_1])])
368         priori_knowledge.add_forbidden_edges([(all_vars[var_1], all_vars[var_2])])
369
370 # Check the priori knowledge
371 print('\nLLM knowledge vs true DAG')
372 priori_dag = np.clip(priori_knowledge.matrix, 0, 1)
373 print(f'\nChecking if priori graph is a DAG: {check_if_dag(priori_dag)}')
374
375 GraphDAG(
376     est_dag=priori_dag,
377     true_dag=true_dag)
378 plt.show()
```

```
379
380  # Now, run the PC algorithm with priori knowledge
381  print('\nRunning PC with Priori Knowledge')
382
383  # Instantiate the model with expert knowledge
384  pc_priori = PC(
385      priori_knowledge=priori_knowledge,
386      variant='stable'
387  )
388
389  # Learn using the dataset
390  pc_priori.learn(data_subset)
391
392  # Visualize
393  GraphDAG(
394      est_dag=pc_priori.causal_matrix,
395      true_dag=true_dag)
396  plt.show()
397
398  # Compute metrics
399  metrics_priori = MetricsDAG(
400      B_est=pc_priori.causal_matrix,
401      B_true=true_dag)
402  print(metrics_priori.metrics)
```

## B.2    Causal Model

```
1   import dowhy
2   from dowhy import CausalModel
3
4   import numpy as np
5   import pandas as pd
6   import graphviz
7   import networkx as nx
8
9   np.set_printoptions(precision=3, suppress=True)
10  np.random.seed(0)
11
12  file_path = r"C:\Users\nstep\TSU\SeniorProject\df_selected1.csv"
13  nashvilleDF = pd.read_csv(file_path)
14
15  def make_graph(adjacency_matrix, labels=None):
16      idx = np.abs(adjacency_matrix) > 0.01
17      dirs = np.where(idx)
18      d = graphviz.Digraph(engine='dot')
19      names = labels if labels else [f'x{i}' for i in range(len(adjacency_matrix))]
20      for name in names:
21          d.node(name)
22      for to, from_, coef in zip(dirs[0], dirs[1], adjacency_matrix[idx]):
23          d.edge(names[from_], names[to], label=str(coef))
24      return d
25
26  def str_to_dot(string):
27      '''
```

```
28      Converts input string from graphviz library to valid DOT graph format.
29      '''
30      graph = string.strip().replace('\n', ';').replace('\t','')
31      graph = graph[:9] + graph[10:-2] + graph[-1] # Removing unnecessary characters from string
32      return graph
33
34  print(nashvilleDF.shape)
35
36  nashvilleDF.head()
37
38
39  # Drop rows where any of the specified columns has a value of 1
40  columns_to_check = ['prop_Entire condo', 'prop_Entire guest suite', 'prop_Entire guesthouse',
41                      'prop_Entire rental unit', 'prop_Entire townhouse', 'prop_Hotel',
42                      'prop_Private room', 'free parking', 'room_Private room']
43
44  nashvilleDF = nashvilleDF[~nashvilleDF[columns_to_check].eq(1).any(axis=1)]
45
46  # Drop 'id' and the specified columns
47  columns_to_drop = ['id','prop_Entire home', 'room_Entire home/apt',] + columns_to_check
48  nashvilleDF = nashvilleDF.drop(columns=columns_to_drop)
49  nashvilleDF.head()
50
51  import graphviz
52
53  # Given DOT representation
54  causal_graph_updated = f"""
55  digraph {{
56      grill -> price;
57      pool -> price;
58      "hot tub" -> price;
59      "hot tub" -> bathrooms;
60      pool -> bathrooms;
61      "private entrance" -> price;
62      "private entrance" -> "resort access";
63      fireplace -> price;
64
65      neighbourhood_cleansed_num -> price;
66      neighbourhood_cleansed_num -> review_scores_location;
67      review_scores_location -> review_scores_rating;
68      review_scores_location -> reviews_per_month;
69      review_scores_rating -> reviews_per_month;
70      reviews_per_month -> price;
71      host_is_superhost -> reviews_per_month;
72      host_is_superhost -> price;
73      host_is_superhost -> review_scores_rating;
74      host_acceptance_rate -> reviews_per_month;
75      host_acceptance_rate -> price;
76      beds -> bedrooms;
77      bedrooms -> accommodates;
78      bedrooms -> bathrooms;
79      accommodates -> price;
80      "resort access" -> pool;
81      "resort access" -> "hot tub";
82      pool -> host_acceptance_rate;
83      "hot tub" -> host_acceptance_rate;
```

```
 84     bathrooms -> price;
 85     minimum_minimum_nights -> price;
 86     minimum_maximum_nights -> price;
 87
 88     "resort access" -> price;
 89 }}
 90 """
 91
 92
 93 model=CausalModel(
 94         data = nashvilleDF,
 95         treatment='accommodates',
 96         outcome='price',
 97         graph=causal_graph_updated)
 98
 99 model.view_model()
100 from IPython.display import Image, display
101 display(Image(filename="causal_model.png"))
102
103
104 # Identification
105 identified_estimand = model.identify_effect(proceed_when_unidentifiable=True)
106 print(identified_estimand)
107
108
109
110 from cdt.causality.graph import LiNGAM, PC, GES, SAM,CAM
111 import graphviz
112 from sklearn.preprocessing import StandardScaler
113
114
115 import os
116 os.environ['R_HOME'] = 'C:\\Program Files\\R\\R-4.3.1'
117 os.environ["PATH"] += os.pathsep + 'C:\\Program Files\\Graphviz\\bin'
118
119 scaler = StandardScaler()
120 nashvilleDF_scaled = pd.DataFrame(scaler.fit_transform(nashvilleDF), columns=nashvilleDF.columns)
121 from sklearn.decomposition import PCA
122
123 pca = PCA(n_components=min(nashvilleDF_scaled.shape) - 1)  # -1 to ensure the matrix is non-singular
124 nashvilleDF_pca = pca.fit_transform(nashvilleDF_scaled)
125
126
127
128
129 graphs = {}
130 labels = [f'{col}' for i, col in enumerate(nashvilleDF.columns)]
131 functions = {
132     'LiNGAM' : LiNGAM,
133     'PC': PC,
134     'GES': GES,
135     'SAM': SAM,
136 }
137
138
139 for method, lib in functions.items():
```

```
140    obj = lib()
141    output = obj.predict(nashvilleDF_scaled)
142    adj_matrix = nx.to_numpy_array(output)
143    adj_matrix = np.asarray(adj_matrix)
144    graph_dot = make_graph(adj_matrix, labels)
145    graphs[method] = graph_dot
146
147 # Visualize graphs
148 for method, graph in graphs.items():
149     print("Method : %s" % (method))
150     display(graph)
151
152
153
154 for method, graph in graphs.items():
155         if method != "LiNGAM":
156             continue
157         print('\n*****************************************************************************\n')
158         print("Causal Discovery Method : %s"%(method))
159
160         # Obtain valid dot format
161         graph_dot = str_to_dot(graph.source)
162
163         # Define Causal Model
164         model=CausalModel(
165                 data = nashvilleDF,
166                 treatment='accommodates',
167                 outcome='price',
168                 graph=graph_dot)
169
170         # Identification
171         identified_estimand = model.identify_effect(proceed_when_unidentifiable=True)
172         print(identified_estimand)
173
174         # Estimation
175         estimate = model.estimate_effect(identified_estimand,
176                                          method_name="backdoor.linear_regression",
177                                          control_value=0,
178                                          treatment_value=1,
179                                          confidence_intervals=True,
180                                          test_significance=True)
181         print("Causal Estimate is " + str(estimate.value))
182
183
184
185
186
187 import graphviz
188
189 # Given DOT representation
190 causal_graph_updated = f"""
191 digraph {{
192     grill -> price;
193     private_entrance -> price;
194     pool -> price;
195     pool -> grill;
```

```
196     "hot tub" -> price;
197     "hot tub" -> bathrooms;
198     pool -> bathrooms;
199     "private entrance" -> price;
200     "private entrance" -> "resort access";
201     fireplace -> price;
202     prop_Entire_Home -> price;
203
204     neighbourhood_cleansed_num -> price;
205     neighbourhood_cleansed_num -> review_scores_location;
206
207     review_scores_location -> reviews_per_month;
208     review_scores_rating -> reviews_per_month;
209     reviews_per_month -> price;
210     host_is_superhost -> reviews_per_month;
211     host_is_superhost -> price;
212     host_is_superhost -> review_scores_rating;
213     host_acceptance_rate -> reviews_per_month;
214     host_acceptance_rate -> price;
215     beds -> bedrooms;
216     bedrooms -> accommodates;
217     accommodates -> bathrooms;
218     accommodates -> price;
219     "resort access" -> pool;
220     "resort access" -> "hot tub";
221     bathrooms -> price;
222
223     minimum_minimum_nights -> price;
224     minimum_maximum_nights -> price;
225
226     "resort access" -> price;
227 }}
228 """
229
230
231
232 #visualize the updated causal graph
233 graph = graphviz.Source(causal_graph_updated)
234 graph # display the graph
235
236
237
238
239
240
241 # Define Causal Model
242 model=CausalModel(
243         data = nashvilleDF,
244         treatment='accommodates',
245         outcome='price',
246         graph=causal_graph_updated)
247
248 # Identification
249 identified_estimand = model.identify_effect(proceed_when_unidentifiable=True)
250 print(identified_estimand)
251
```

```python
252  # Estimation
253  estimate = model.estimate_effect(identified_estimand,
254                                   method_name="backdoor.linear_regression",
255                                   control_value=0,
256                                   treatment_value=1,
257                                   confidence_intervals=True,
258                                   test_significance=True)
259  print("Causal Estimate is " + str(estimate.value))
260
261
262
263  # Refutation
264  # Add a random common cause variable
265  refutation = model.refute_estimate(identified_estimand, estimate,
266                                     method_name="data_subset_refuter")
267  print(refutation)
268
269
270  # Refutation
271  # Add a random common cause variable
272  refutation = model.refute_estimate(identified_estimand, estimate,
273                                     method_name="random_common_cause")
274  print(refutation)
275
276  from dowhy.causal_estimators import CausalEstimator
277  from dowhy.causal_refuters.data_subset_refuter import DataSubsetRefuter
278  from upload_data_page import upload_data
279  from feature_engineering import feature_engineering, display_data_preview, display_handle_missing_values,
         display_process_currency_percentage
280  from feature_engineering import display_drop_columns, display_data_transformation, display_encode_categorical,
         display_time_series_features, display_convert_to_datetime
281  from explore_data import explore_data, display_boxplot, display_binary_distribution, feature_comparison_graph_page
282  from regression import evaluate_model_page, display_model_performance_comparison, prepare_data, create_models,
         fit_models, evaluate_models, plot_model_performance
283  from regression import display_select_target_features_and_train, display_feature_importance,
         display_prediction_vs_actual, display_residuals_plot
284  from regression import display_correlation_heatmap, evaluate_model_page
285  from advance_data_analysis import advanced_data_analysis, perform_classification, perform_clustering,
         perform_dimensionality_reduction
286  from time_series_analysis import time_series_analysis, visualize_time_series_data,display_acf_pacf, fit_arima_model
287  import streamlit as st
288  import pandas as pd
289  import numpy as np
290  from sklearn.preprocessing import StandardScaler, MinMaxScaler, LabelEncoder
291  from sklearn.model_selection import train_test_split
292  from sklearn.metrics import r2_score, mean_squared_error, mean_absolute_error
293  from sklearn.ensemble import GradientBoostingRegressor, RandomForestClassifier, RandomForestRegressor,
         HistGradientBoostingRegressor
294  from sklearn.linear_model import LinearRegression
295  from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
296  import xgboost as xgb
297  import seaborn as sns
298  import matplotlib.pyplot as plt
299  import streamlit as st
300  import pandas as pd
301  from sklearn.model_selection import train_test_split
```

```python
302  from sklearn.ensemble import RandomForestClassifier
303  from sklearn.cluster import KMeans
304  from sklearn.decomposition import PCA
305  from sklearn.feature_selection import SelectKBest, f_classif
306  import plotly.express as px
307  from sklearn.preprocessing import LabelEncoder
308  import seaborn as sns
309  import base64
310  from prophet import Prophet
311  import os
312  from dowhy import CausalModel
313  import re
314  import graphviz
315  from graphviz import Digraph
316  from io import BytesIO
317
318  def extract_relationships_from_dot(dot_representation):
319      """Extract relationships from DOT format."""
320      relationships = []
321
322      # Split the string into lines and filter out non-edge lines
323      lines = dot_representation.split("\n")
324      edge_lines = [line.strip() for line in lines if "->" in line]
325
326      for edge_line in edge_lines:
327          # Extract cause and effect from the edge line
328          cause, effect = edge_line.split("->")
329          relationships.append((cause.strip(), effect.strip()))
330
331      return relationships
332
333  def display_relationships_definition():
334      st.subheader("Define Causal Relationships")
335
336      # Ensure data is uploaded
337      if 'data' not in st.session_state or st.session_state.data is None:
338          st.warning("Please upload data first.")
339          return
340
341      columns = list(st.session_state.data.columns)
342
343      # Dropdowns to select cause and effect columns
344      cause_column = st.selectbox("Select Cause Column", columns)
345      effect_column = st.selectbox("Select Effect Column", columns)
346
347      # Add relationship button
348      if st.button("Add Relationship"):
349          relationship = (cause_column, effect_column)
350          if "relationships" not in st.session_state:
351              st.session_state.relationships = []
352          if relationship not in st.session_state.relationships:
353              st.session_state.relationships.append(relationship)
354              st.success(f"Added relationship: {cause_column} -> {effect_column}")
355
356      # Display existing relationships and allow removal
357      if "relationships" in st.session_state and st.session_state.relationships:
```

```
358         st.write("Defined Relationships:")
359
360         st.write("Select Relationship to Remove:")
361         # Create a dropdown with all relationships formatted as "cause -> effect"
362         relationship_options = [f"{cause} -> {effect}" for cause, effect in st.session_state.relationships]
363         selected_relationship_str = st.selectbox("Select Relationship:", relationship_options)
364
365         # Extract cause and effect from the selected string
366         selected_cause, selected_effect = selected_relationship_str.split(" -> ")
367
368         if st.button(f"Remove {selected_cause} -> {selected_effect}"):
369             st.session_state.relationships.remove((selected_cause, selected_effect))
370             st.success(f"Removed relationship: {selected_cause} -> {selected_effect}")
371
372     # Upload a .dot file and load the graph
373     uploaded_file = st.file_uploader("Upload a .dot file to load a causal graph", type="dot")
374     if uploaded_file:
375         uploaded_graph = uploaded_file.read().decode()
376
377         # Extract relationships from the uploaded graph
378         uploaded_relationships = extract_relationships_from_dot(uploaded_graph)
379
380         # Add the extracted relationships to st.session_state.relationships
381         if "relationships" not in st.session_state:
382             st.session_state.relationships = []
383
384         for relationship in uploaded_relationships:
385             if relationship not in st.session_state.relationships:
386                 st.session_state.relationships.append(relationship)
387
388         st.success(f"Loaded {len(uploaded_relationships)} relationships from the provided .dot file.")
389
390     # Generate causal graph button
391     if st.button("Generate Causal Graph"):
392         if "relationships" not in st.session_state or not st.session_state.relationships:
393             st.warning("Please define at least one relationship before generating the graph.")
394         else:
395             # Generate graph using Graphviz
396             dot = Digraph()
397             for cause, effect in st.session_state.relationships:
398                 dot.edge(cause, effect)
399
400             # Convert dot to string and save in session state
401             st.session_state.dot_representation = dot.source
402             st.session_state.generated_graph = True
403
404             # Display the graph
405             st.graphviz_chart(dot)
406
407             # Provide a download link for the graph
408             # Save the graph representation as a temporary .dot file
409             temp_filename = "causal_graph.dot"
410             with open(temp_filename, "w") as f:
411                 f.write(dot.source)
412
413             st.markdown(generate_download_link(temp_filename, "Download causal graph (.dot)"), unsafe_allow_html=True)
```

```
414
415
416  def display_causal_model_creation():
417      """Sub-task for creating the causal model based on the defined graph."""
418
419      columns = list(st.session_state.data.columns)
420
421      # Check if the causal graph has been generated
422      if not st.session_state.get("generated_graph", False):
423          st.warning("Please generate the causal graph first.")
424          return
425
426      # Ensure that the dot_representation is not empty
427      dot_representation = st.session_state.get("dot_representation", "")
428      if not dot_representation:
429          st.warning("Please generate or upload a causal graph first.")
430          return
431
432      # Causal model creation
433      treatment = st.selectbox("Select Treatment (cause) Variable", columns)
434      outcome = st.selectbox("Select Outcome (effect) Variable", columns)
435      st.write("""
436      The treatment variable is what you believe to be the cause in your causal relationship,
437      and the outcome variable is the effect you are studying.
438      """)
439
440      if st.button("Create and Estimate Causal Model"):
441          # Define Causal Model
442          model = CausalModel(
443              data=st.session_state.data,
444              treatment=treatment,
445              outcome=outcome,
446              graph=st.session_state.get("dot_representation", "")
447          )
448
449
450          # Identification
451          identified_estimand = model.identify_effect(proceed_when_unidentifiable=True)
452          st.session_state.identified_estimand = identified_estimand
453          st.write("Identified estimand:", identified_estimand)
454
455          # Estimation
456          estimate = model.estimate_effect(identified_estimand,
457                                           method_name="backdoor.linear_regression",
458                                           control_value=0,
459                                           treatment_value=1,
460                                           confidence_intervals=True,
461                                           test_significance=True)
462          st.write("Causal Estimate:", estimate.value)
463
464          st.session_state.causal_model = model
465          st.session_state.estimate = estimate
466          st.success("Causal model created and estimated successfully!")
467
468
469
```

```
470
471 def generate_download_link(filename, download_text):
472     """Generate a download link for a given file and link text."""
473     with open(filename, "rb") as f:
474         file_data = f.read()
475     b64 = base64.b64encode(file_data).decode()
476     href = f'<a href="data:application/octet-stream;base64,{b64}" download="{filename}">{download_text}</a>'
477     return href
478
479
480
481
482
483
484
485
486 def display_refutation_tests():
487     """Sub-task for running refutation tests."""
488     st.subheader("Refutation Tests")
489
490     # Ensure that the causal model is created
491     if "causal_model" not in st.session_state:
492         st.warning("Please create a causal model first.")
493         return
494
495     # Refutation methods
496     methods = ["Data Subset Refuter"]
497     chosen_method = st.selectbox("Choose a Refutation Method", methods)
498
499     # Customize parameters based on the chosen method
500     if chosen_method == "Data Subset Refuter":
501         subset_fraction = st.slider("Choose a fraction of data to keep", 0.1, 1.0, 0.5)
502
503     # Run refutation
504     if st.button("Run Refutation"):
505         if "identified_estimand" not in st.session_state or "estimate" not in st.session_state:
506             st.warning("Please create and estimate the causal model first.")
507             return
508
509         refuter = DataSubsetRefuter(
510             data=st.session_state.data,
511             causal_model=st.session_state.causal_model,
512             identified_estimand=st.session_state.identified_estimand,
513             estimate=st.session_state.estimate,
514             subset_fraction=subset_fraction
515         )
516         results = refuter.refute_estimate()
517         st.write("Refutation Results:", results)
518
519         # Extract p_value from the results
520         # Assuming it's in the results string, you might need to adjust the extraction method
521         p_value_str = re.search(r'p value:(\d+\.\d+)', str(results))
522         if p_value_str:
523             p_value = float(p_value_str.group(1))
524         else:
525             p_value = None
```

```
526
527        # Interpretation based on p-value and difference in effects
528        original_effect = st.session_state.estimate.value
529        new_effect = results.new_effect
530
531        if p_value and p_value > 0.05 and abs(original_effect - new_effect) < 0.05 * abs(original_effect):  # Assuming a
      5% relative difference threshold for "close"
532            st.write("Interpretation: The original causal estimate is consistent and robust, even when using a subset of
      the data.")
533        elif p_value and p_value <= 0.05:
534            st.write("Interpretation: The original causal estimate may not be reliable, as it changes significantly with
      a subset of the data.")
535
536
537
538
539
540
541 def causality_page():
542    st.header("Causality Analysis")
543
544    # Ensure data is uploaded
545    if 'data' not in st.session_state or st.session_state.data is None:
546        st.warning("Please upload data first.")
547        return
548
549    # Moved task selection to sidebar
550    task = st.sidebar.radio("Choose a Causality Sub-task", ["Define Relationships", "Create Causal Model", "Run
      Refutation Tests"])
551
552    if task == "Define Relationships":
553        display_relationships_definition()
554    elif task == "Create Causal Model":
555        display_causal_model_creation()
556    elif task == "Run Refutation Tests":
557        display_refutation_tests()
```

# WEB APPLICATION

## C.1 Prediction Page

```html
1  <!DOCTYPE html>
2  <html lang="en">
3
4  <head>
5      <meta charset="UTF-8">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7
8      <!-- Allows for the styling of the web page -->
9      <link rel="stylesheet" type="text/css" href="{{ url_for('static', filename='Styles/style.css') }}">
10
11     <title>Predictor</title>
12 </head>
13
14 <body>
15     <!-- The tool bar on the web app -->
16     <div class="topnav">
17         <a class="active" href="{{ url_for('predictor') }}">Predictor</a>
18     </div>
19
20     <h1 >Price Prediction</h1>
21
22     <form method="POST" action="{{ url_for('predictor') }}">
23         {{ form.hidden_tag() }}
24         <table>
25             <th>
26                 {{ form.accommodates.label }}<br>{{ form.accommodates(class="input-field") }}<br>
27                 {% if form.accommodates.errors %}
28                     <div class="error">
29                     {% for error in form.accommodates.errors %}
30                         <p>Invalid character. Please enter number.</p>
31                     {% endfor %}
32                     </div>
33                 {% endif %}
34                 {{ form.bathrooms.label }}<br>{{ form.bathrooms(class="input-field") }}<br>
35                 {% if form.bathrooms.errors %}
36                     <div class="error">
37                     {% for error in form.bathrooms.errors %}
38                         <p>Invalid character. Please enter number.</p>
39                     {% endfor %}
40                     </div>
41                 {% endif %}
42                 {{ form.nights_staying.label }}<br>{{ form.nights_staying(class="input-field") }}<br>
```

```
43                  {% if form.nights_staying.errors %}
44                      <div class="error">
45                      {% for error in form.nights_staying.errors %}
46                          <p>Invalid character. Please enter number.</p>
47                      {% endfor %}
48                      </div>
49                  {% endif %}
50                  {{ form.bedrooms.label }}<br>{{ form.bedrooms(class="input-field") }}<br>
51                  {% if form.bedrooms.errors %}
52                      <div class="error">
53                      {% for error in form.bedrooms.errors %}
54                          <p>Invalid character. Please enter number.</p>
55                      {% endfor %}
56                      </div>
57                  {% endif %}
58                  {{ form.property_type.label }}<br>{{ form.property_type(class="input-field") }}<br>
59              </th>
60
61              <th>
62                  {{ form.beds.label }}<br>{{ form.beds(class="input-field") }}<br>
63                  {% if form.beds.errors %}
64                      <div class="error">
65                      {% for error in form.beds.errors %}
66                          <p>Invalid character. Please enter number.</p>
67                      {% endfor %}
68                      </div>
69                  {% endif %}
70                  {{ form.review_scores_rating.label }}<br>{{ form.review_scores_rating(class="input-field") }}<br>
71                  {% if form.review_scores_rating.errors %}
72                      <div class="error">
73                      {% for error in form.review_scores_rating.errors %}
74                          <p>Invalid character. Please enter number.</p>
75                      {% endfor %}
76                      </div>
77                  {% endif %}
78                  {{ form.review_scores_location.label }}<br>{{ form.review_scores_location(class="input-field") }}<br>
79                  {% if form.review_scores_location.errors %}
80                      <div class="error">
81                      {% for error in form.review_scores_location.errors %}
82                          <p>Invalid character. Please enter number.</p>
83                      {% endfor %}
84                      </div>
85                  {% endif %}
86                  {{ form.host_acceptance_rate.label }}<br>{{ form.host_acceptance_rate(class="input-field") }}<br>
87                  {% if form.host_acceptance_rate.errors %}
88                      <div class="error">
89                      {% for error in form.host_acceptance_rate.errors %}
90                          <p>Invalid character. Please enter number.</p>
91                      {% endfor %}
92                      </div>
93                  {% endif %}
94                  {{ form.room_type.label }}<br>{{ form.room_type(class="input-field") }}<br>
95              </th>
96
97              <th>
98                  {{ form.fireplace.label }}<br>{{ form.fireplace() }}<br>
```

93

```
99                {{ form.hot_tub.label }}<br>{{ form.hot_tub() }}<br>
100               {{ form.grill.label }}<br>{{ form.grill() }}<br>
101               {{ form.pool.label }}<br>{{ form.pool() }}<br>
102               {{ form.private_entrance.label }}<br>{{ form.private_entrance() }}<br>
103               {{ form.free_parking.label }}<br>{{ form.free_parking() }}<br>
104               {{ form.resort_access.label }}<br>{{ form.resort_access() }}<br>
105               {{ form.neighbourhood_cleansed_num.label }}<br>{{ form.neighbourhood_cleansed_num(class="input-field")
      }}<br>
106           </th>
107       </table>
108       {{ form.submit(class="submit-button") }}<br>
109
110       <center>
111           <img
       src="https://filetransfer.nashville.gov/portals/0/sitecontent/Planning/images/FrontPage/DistrictsAndCommunities600.jpg"
       alt="Nashville District Map">
112       </center>
113   </form>
114
115   {% if estimated_price %}
116       <div class="results-container">
117           <div class="result">
118               Predicted Price: ${{ estimated_price }}
119           </div>
120       </div>
121   {% endif %}
122 </body>
123
124 </html>
```

## C.2      Price Page

```
1  <!DOCTYPE html>
2  <html lang="en">
3
4  <head>
5      <meta charset="UTF-8">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7
8      <!-- Allows for the styling of the web page -->
9      <link rel="stylesheet" type="text/css" href="{{ url_for('static', filename='Styles/style.css') }}">
10
11     <title>Airbnb Price Result</title>
12 </head>
13 <body>
14     <!-- The tool bar on the web app -->
15     <div class="topnav">
16         <a href="{{ url_for('predictor') }}">Predictor</a>
17         <a class="active" href="{{ url_for('price') }}">Price</a>
18     </div>
19
20     <div class="price-container">
21         <div class="price-content">
22             <p>Predicted Price: ${{ estimated_price }}</p>
```

```
23            <p>MAE: ${{ mae }} (This indicates the average deviation of the predicted price from actual prices. A lower
      value means the prediction is more accurate.)</p>
24        </div>
25      </div>
26
27      <hr>
28      <h2 style="text-align: center;">Recommended Airbnb Listings:</h2>
29      <div style="text-align: center; margin-top: 20px;">
30          {% for prop in recommended_properties.iterrows() %}
31              <div>
32                  <a href="https://www.airbnb.com/rooms/{{ prop[1]['id']|int }}">Property ID: {{ prop[1]['id']|int }}</a>
33              </div>
34          {% endfor %}
35      </div>
36  </body>
37  </html>
```

## C.3  Style Sheet

```
1  body {
2    height: 100%;
3    background-color: #ffffff;
4    color: rgb(0, 0, 0);
5    font-family:Cambria, Cochin, Georgia, Times, 'Times New Roman', serif;
6  }
7
8   /* Adds style to the input fields */
9  input[type=text] {
10    width: 50%;
11    padding: 12px 20px;
12    margin: 8px 0;
13    box-sizing: border-box;
14    border: 2px solid black;
15    border-radius: 4px;
16    background-color: rgba(128, 128, 128, 0.281);
17  }
18
19  /* Adds style to the buttons */
20  button {
21    background-color: #6e6e6e51;
22    border: 2px solid black;
23    color: rgb(250, 248, 248);
24    padding: 15px 32px;
25    text-align: center;
26    text-decoration: none;
27    display: inline-block;
28    font-size: 12px;
29  }
30
31  table {
32    width: 100%;
33    height: 400px;
34    border-collapse: collapse;
35  }
```

```
36  th, td {
37    padding: 15px;
38    text-align: left;
39    border-bottom: 1px solid #ddd;
40    text-align: center;
41  }
42
43  /* The error statement will show as Red */
44  .error {
45    color: red;
46  }
47
48  /* Adds style to the drop down menus */
49  select {
50    appearance: none;
51    background-color: rgba(128, 128, 128, 0.281);
52    border: 2px solid black;
53    border-radius: 4px;
54    outline: black;
55    padding: 12px 20px;
56    margin: 8px 0;
57    width: 50%;
58    font-family: inherit;
59    font-size: inherit;
60    color: rgb(0, 0, 0);
61  }
62
63  /* Add a black background color to the top navigation */
64  .topnav {
65    margin: 0;
66    background-color: transparent;
67    overflow: hidden;
68  }
69
70  /* Style the links inside the navigation bar */
71  .topnav a {
72    float: left;
73    color: #000000;
74    text-align: center;
75    padding: 14px 16px;
76    text-decoration: none;
77    font-size: 17px;
78  }
79
80  /* Change the color of links on hover */
81  .topnav a:hover {
82    background-color: #ddd;
83    color: black;
84  }
85
86  /* Add a color to the active/current link */
87  .topnav a.active {
88    background-color: #0d06f3;
89    color: white;
90  }
91
```

```
 92   .column {
 93     float: left;
 94     width: 50%;
 95   }
 96
 97   /* Clear floats after the columns */
 98   .row:after {
 99     content: "";
100     display: table;
101     clear: both;
102   }
103
104   form {
105     text-align: center;
106   }
107
108   /* Center results on the page */
109   .results-container {
110     display: flex;
111     flex-direction: column;
112     align-items: center;
113     justify-content: center;
114     height: calc(100vh - 50px); /* Adjust 50px based on your actual navbar height */
115   }
116
117   .result {
118     font-size: 24px;
119     margin-bottom: 20px;
120   }
121
122   .price-container {
123     text-align: center; /* This will center the content horizontally */
124     padding-top: 5rem; /* Adjust this value to increase/decrease space on top */
125     padding-bottom: 1rem; /* Adjust this value to increase/decrease space at the bottom */
126   }
127
128   .price-content {
129     padding-bottom: 3rem; /* Adjust this value to increase/decrease space below the MAE before the line */
130   }
131
132
133   hr {
134     margin: 0 auto; /* This centers the horizontal line */
135     width: 80%; /* Adjusts the width of the line */
136   }
137
138
139   h2 {
140     margin-top: 3rem; /* Adjust this value to increase/decrease space above the title */
141     margin-bottom: 1rem; /* Adjust this value to increase/decrease space below the title */
142   }
```

## C.4 Back-end

```
1   import pandas as pd
2   from flask import Flask, render_template, request
3   from flask_wtf import FlaskForm
4   from wtforms import FloatField, SubmitField, BooleanField, SelectField
5   from wtforms.validators import DataRequired, InputRequired, NumberRange
6   from joblib import dump, load
7   from sklearn.ensemble import HistGradientBoostingRegressor
8   from sklearn.metrics import mean_absolute_error
9   import warnings
10  from sklearn.model_selection import train_test_split
11  warnings.filterwarnings(action='ignore', category=UserWarning, module='sklearn')
12
13  app = Flask(__name__, static_folder='Static')
14  app.config['SECRET_KEY'] = 'wdazD1dRmBGVwVSi'
15
16  # Load or train the model
17  try:
18      trained_model = load('trained_model.joblib')
19      mean_features = nashvilleDF.drop(['price', 'id'], axis=1).mean().to_dict()
20  except:
21      nashvilleDF = pd.read_csv('nashvilleDF.csv')
22      mean_features = nashvilleDF.mean().to_dict()
23      X = nashvilleDF.drop(['price', 'id'], axis=1)
24      y = nashvilleDF['price']
25      trained_model = HistGradientBoostingRegressor()
26      trained_model.fit(X, y)
27      dump(trained_model, 'trained_model.joblib')
28
29
30
31  X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
32
33  trained_model.fit(X_train, y_train)
34  y_pred = trained_model.predict(X_test)
35  mae = mean_absolute_error(y_test, y_pred)
36
37
38  # Prediction Page Forms
39  class PredictionForm(FlaskForm):
40      accommodates = FloatField('How many people do you want to stay?', validators=[DataRequired(), NumberRange(min=1)])
41      bathrooms = FloatField('Enter the number of bathrooms (e.g. 1.5)', validators=[DataRequired(), NumberRange(min=0)])
42      nights_staying = FloatField('How many nights are you staying?', validators=[DataRequired(), NumberRange(min=1)])
43      grill = BooleanField('Do you want a grill?')
44      bedrooms = FloatField('Enter the number of bedrooms', validators=[DataRequired(), NumberRange(min=0)])
45      fireplace = BooleanField('Do you want a fireplace?')
46      hot_tub = BooleanField('Do you want a hot tub?')
47      private_entrance = BooleanField('Do you want a private entrance?')
48      free_parking = BooleanField('Do you want free parking?')
49      review_scores_rating = FloatField('Enter desired review score (1-5 or percentage)', validators=[DataRequired(),
         NumberRange(min=0, max=100)])
50      review_scores_location = FloatField('Enter desired location review score (1-5 or percentage)',
         validators=[DataRequired(), NumberRange(min=0, max=100)])
51      resort_access = BooleanField('Do you want resort access?')
52      beds = FloatField('Enter number of beds', validators=[DataRequired(), NumberRange(min=1)])
53      pool = BooleanField('Do you want a pool?')
54      host_acceptance_rate = FloatField('Enter desired host acceptance rate (0-100)', validators=[DataRequired(),
```

```
     NumberRange(min=0, max=100)])
55    neighbourhood_cleansed_num = SelectField('Choose Neighbourhood (Chart below)', choices=[
56        (1, 'District 1'), (2, 'District 2'), (3, 'District 3'), (4, 'District 4'), (5, 'District 5'), (6, 'District 6'),
      (7, 'District 7'), (8, 'District 8'),
57        (9, 'District 9'), (10, 'District 10'), (11, 'District 11'), (12, 'District 12'), (13, 'District 13'), (14,
      'District 14'), (15, 'District 15'), (16, 'District 16'),
58        (17, 'District 17'), (18, 'District 18'), (19, 'District 19'), (20, 'District 20'), (21, 'District 21'), (22,
      'District 22'), (23, 'District 23'), (24, 'District 24'),
59        (25, 'District 25'), (26, 'District 26'), (27, 'District 27'), (28, 'District 28'), (29, 'District 29'), (30,
      'District 30'), (31, 'District 31'), (32, 'District 32'),
60        (33, 'District 33'), (34, 'District 34'), (35, 'District 35')
61    ])
62    property_type = SelectField('Choose Property Type', choices=[
63        ('prop_Entire condo', 'Entire Condo'),
64        ('prop_Entire guest suite', 'Entire Guest Suite'),
65        ('prop_Entire guesthouse', 'Entire Guesthouse'),
66        ('prop_Entire home', 'Entire Home'),
67        ('prop_Entire rental unit', 'Entire Rental Unit'),
68        ('prop_Entire townhouse', 'Entire Townhouse'),
69        ('prop_Hotel', 'Hotel'),
70        ('prop_Private room', 'Private Room')
71    ])
72    room_type = SelectField('Choose Room Type', choices=[
73    ('room_Entire home/apt', 'Entire home/apt'),
74    ('room_Private room', 'Private Room')
75 ])
76

77

78    submit = SubmitField('Predict Price')

79

80 # Route to the Prediction Model Page
81 @app.route('/')
82 @app.route('/predictor', methods=['GET', 'POST'])
83 def predictor():
84    form = PredictionForm()
85    estimated_price = None

86

87    if form.validate_on_submit():
88        # Creates the user_data dictionary (extends form.data)
89        user_data = {key: value for key, value in form.data.items() if key != 'nights_staying' and key != 'csrf_token'}

90

91        # Convert specific fields to appropriate values
92        for key, value in user_data.items():
93            if value == 'Yes':
94                user_data[key] = 1
95            elif value == 'No':
96                user_data[key] = 0
97            elif isinstance(value, str) and "%" in value:
98                user_data[key] = float(value.strip('%')) / 100

99

100        # Handle one-hot encoding for property_type
101        property_types = [
102            'prop_Entire condo', 'prop_Entire guest suite', 'prop_Entire guesthouse',
103            'prop_Entire home', 'prop_Entire rental unit', 'prop_Entire townhouse',
104            'prop_Hotel', 'prop_Private room'
105        ]
```

```python
106            selected_property_type = user_data.pop('property_type')
107            for prop_type in property_types:
108                user_data[prop_type] = 1 if prop_type == selected_property_type else 0
109
110            # Ensure the Amenities are correctly encoded
111            if 'hot_tub' in user_data:
112                user_data['hot tub'] = user_data.pop('hot_tub')
113            if 'grill' in user_data:
114                user_data['grill'] = user_data.pop('grill')
115            if 'private_entrance' in user_data:
116                user_data['private entrance'] = user_data.pop('private_entrance')
117            if 'free_parking' in user_data:
118                user_data['free parking'] = user_data.pop('free_parking')
119            if 'resort_access' in user_data:
120                user_data['resort access'] = user_data.pop('resort_access')
121            if 'pool' in user_data:
122                user_data['pool'] = user_data.pop('pool')
123
124            # Populate missing features with their average values
125            prepared_data = mean_features.copy()
126            prepared_data.update(user_data)
127
128            if 'id' in prepared_data:
129                del prepared_data['id']
130            if 'price' in prepared_data:
131                del prepared_data['price']
132            if 'submit' in prepared_data:
133                del prepared_data['submit']
134
135            # Now, let's filter the nashvilleDF dataframe based on the user's criteria:
136
137            # Start with the entire dataset and filter it step by step
138            filtered_properties = nashvilleDF.copy()
139
140            # 1. Accommodates ``
141            filtered_properties = filtered_properties[filtered_properties['accommodates'] >= user_data['accommodates']]
142
143            # 2. Bathrooms
144            filtered_properties = filtered_properties[filtered_properties['bathrooms'] >= user_data['bathrooms']]
145
146            # 3. Nights staying
147            # Use nights_staying to filter the properties
148            nights_staying_value = form.data['nights_staying']
149            filtered_properties = filtered_properties[
150                (filtered_properties['minimum_minimum_nights'] <= nights_staying_value) &
151                (filtered_properties['minimum_maximum_nights'] >= nights_staying_value)
152            ]
153
154            # Now, drop the key from the dictionary since it's not needed anymore
155            user_data.pop('nights_staying', None)
156
157
158            # 4. Bedrooms
159            filtered_properties = filtered_properties[filtered_properties['bedrooms'] >= user_data['bedrooms']]
160
161            # Ensuring the amenities correctly encoded
```

```
162
163
164         # 5. Amenities
165         amenities = ['fireplace', 'hot tub', 'grill', 'private entrance', 'free parking', 'resort access', 'pool']
166         for amenity in amenities:
167             if user_data[amenity]:
168                 filtered_properties = filtered_properties[filtered_properties[amenity] == 1]
169
170         # 6. Review scores
171         min_review_score = user_data['review_scores_rating'] - 1
172         filtered_properties = filtered_properties[filtered_properties['review_scores_rating'] > min_review_score]
173
174         # 7. Property type
175         # The filtering for property type is inherently done through the one-hot encoding process
176
177         # 8. Room type
178         # Handle one-hot encoding for room_type
179         room_types = {
180             'room_Entire home/apt': 'Entire home/apt',
181             'room_Private room': 'Private Room'
182         }
183         selected_room_type = user_data.pop('room_type')
184         for encoded_name, original_name in room_types.items():
185             user_data[encoded_name] = 1 if original_name == selected_room_type else 0
186
187         # 9. Location Review Scores
188         min_review_score_location = user_data['review_scores_location'] - 1
189         filtered_properties = filtered_properties[filtered_properties['review_scores_location'] >
    min_review_score_location]
190
191         # 10. Beds
192         filtered_properties = filtered_properties[filtered_properties['beds'] >= user_data['beds']]
193
194         # Neighbourhood Cleansed
195         neighbourhood = [
196             'District 1', 'District 2', 'District 3', 'District 4', 'District 5', 'District 6', 'District 7', 'District
    8',
197             'District 9', 'District 10', 'District 11', 'District 12', 'District 13', 'District 14', 'District 15',
    'District 16',
198             'District 17', 'District 18', 'District 19', 'District 20', 'District 21', 'District 22', 'District 23',
    'District 24',
199             'District 25', 'District 26', 'District 27', 'District 28', 'District 29', 'District 30', 'District 31',
    'District 32',
200             'District 33', 'District 34', 'District 35'
201         ]
202         selected_neighbuorhood = user_data.pop('neighbourhood_cleansed_num')
203         for neighbourhood_num in neighbourhood:
204             user_data[neighbourhood_num] = 1 if neighbourhood_num == selected_neighbuorhood else 0
205
206         # Host Acceptance Rate
207         filtered_properties = filtered_properties[filtered_properties['host_acceptance_rate'] >=
    user_data['host_acceptance_rate']]
208
209         # Compare features between model and user data
210         model_features = set(X.columns)
211         prepared_features = set(prepared_data.keys())
```

```python
212
213          missing_features = model_features - prepared_features
214          extra_features = prepared_features - model_features
215
216          print("Missing features:", missing_features)
217          print("Extra features:", extra_features)
218
219          # Predict the price using the trained model
220          try:
221              ordered_data_values = [prepared_data[feature] for feature in X.columns]
222
223              estimated_price = trained_model.predict([ordered_data_values])[0]
224              estimated_price = round(estimated_price, 2)
225
226              recommended_properties = filtered_properties.sort_values(by='review_scores_rating', ascending=False).head(5)
227
228              return render_template('price.html', estimated_price=estimated_price, mae=round(mae, 2),
        recommended_properties=recommended_properties)
229
230          except Exception as e:
231              print(f"Error during prediction: {str(e)}")
232
233      return render_template('prediction.html', title='Prediction Model', form=form, estimated_price=estimated_price)
234
235  @app.route('/price')
236  def price():
237      estimated_price = request.args.get('estimated_price', 'N/A')
238      mae = request.args.get('mae', 'N/A')
239      return render_template('price.html', estimated_price=estimated_price, mae=mae)
240
241
242  if __name__ == '__main__':
243      app.run(debug=True)
```

# BIBLIOGRAPHY

[1]   Inside Airbnb. "Airbnb - Listings". In: (2021), https://public.opendatasoft.com/explore/dataset/airbnb–listings/information/?disjunctive.host_verificationsdisjunctive.amenitiesdisjunctive.features.

[2]   Peter Bruce and Andrew Bruce. "Practical Statistics for Data Scientists". In: *O'Reilly* (2017).

[3]   D. Mackenzi J. Pearl. "The book of why: The new science of cause and effect". In: *New York, NY: Basic Books* (2020).

[4]   Emre Kiciman and Amit Sharama. "Dowhy: An end-to-end library for causal inference". In: *arXiv.org* 2020.04216 (2020).

[5]   Rik Kraan. "Demystifying decision trees, random forests  gradient boosting". In: *Towards Data Science* (2020).

[6]   A. Molak. "Jane the Discoverer: Enhancing Causal Discovery with Large Language Models (Causal Python)". In: (2023), https://towardsdatascience.com/jane–the–discoverer–enhancing–causal–discovery–with–large–language–models–causal–python–564a63425c93.

[7]   J. Pearl and D. Mackenzie. "The book of why: The new science of cause and effect". In: *New York: Basic Books* (2020).

[8]   V. Raul Perez-Sanchez Raul-Tomas Mora-Garcia Maria-Francisace Cespedes-Lopez. "Housing Price Prediction Using Machine Learning Algorithms in COVID-19 Times". In: *MDPI* (2022).

[9]   VB Staff. "Airbnb VP talks about AI's profound impact on results". In: *Venture Beat* (2017).

[10]  Naveen Venkat. *The Curse of Dimensionality: Inside Out.* `https://www.researchgate.net/publication/327498046_The_Curse_of_Dimensionality_Inside_Out`. Accessed April 1, 2019.