

CS7CS4 – Machine Learning Supplemental Assignment 2023-24

Xin Wang

School of Computer Science and Statistics

Trinity College Dublin

Dublin, Ireland

Email: wangx33@tcd.ie

Abstract—This study aims to predict the ratings of Airbnb listings in Dublin using machine learning models. By examining various features and services of the listings, we seek to create an accurate predictive model that can assist hosts in improving their offerings and provide renters with dependable information. Data from the Inside Airbnb website, including listings and reviews, were processed and analyzed. We employed and evaluated three machine learning models: Logistic Regression, k-Nearest Neighbors (k-NN), and Decision Tree Classifiers. Logistic Regression consistently outperformed the other models, demonstrating its effectiveness in this context. The findings emphasize the importance of feature engineering and model tuning in achieving high predictive accuracy. This research offers both academic insights into the application of machine learning in real-world situations and practical recommendations for enhancing the Airbnb platform.

Index Terms—Airbnb ratings prediction, Machine learning, Logistic Regression, k-Nearest Neighbors (k-NN), Decision Tree Classifier, Feature engineering, Data preprocessing, Model evaluation

I. INTRODUCTION

Airbnb, a web-based service, has significantly transformed how people search for places to stay and travel. It provides a broad range of lodging options, from whole apartments to individual rooms, catering to diverse travelers' needs. Users can make reservations through the platform and rate their stays on various criteria like cleanliness, communication, overall experience, location, value, and accuracy. These ratings impact the listing's ranking and visibility and are essential in influencing prospective renters' choices.

However, Airbnb ad ratings are shaped by various factors, including the particular features and services of each listing. Predicting these ratings is crucial as it aids hosts in enhancing their offerings and gives renters a more dependable foundation for their decisions. Hence, this paper aims to forecast the possible ratings of Airbnb listings in Dublin by examining their features and services.

To achieve this goal, this study uses data from the Inside Airbnb website. The datasets used include the ads dataset and the reviews dataset. The ads dataset contains 75 feature columns with a total of 7,345 records, while the reviews dataset contains 6 feature columns with a total of 230,065 reviews. Through in-depth analysis and modeling of this data, we hope to build an effective predictive model that accurately predicts the ratings of Airbnb ads in Dublin.

The academic significance of this study stems from applying independent machine learning analyses to thoroughly and comprehensively explore the data. By utilizing various

machine learning algorithms, we identified and validated the key factors affecting listing ratings. This approach not only highlights the potential of machine learning for real-world applications but also showcases the effectiveness of data-driven analytics in addressing complex issues.

This study holds both academic value and practical application. By uncovering the key factors that influence ratings, we can provide valuable insights for enhancing the Airbnb platform and offer a better service experience to both hosts and renters. In summary, this research not only presents a new analytical perspective for academia but also offers significant practical references.

II. FEATURE ENGINEERING

The Listings dataset comprises 75 feature columns, 10 of which relate to data scraping details and audit information (such as **listing URL**, **scrape ID**, **last scraped**, **source**, **picture URL**, **host ID**, **host URL**, **host thumbnail URL**, **host picture URL**, and **calendar last scraped**). These columns were excluded to enhance the dataset's clarity and comprehensibility. Similarly, from the Reviews dataset, only the **listing ID** and **comments** were retained for further analysis from the initial 6 features. Detailed preprocessing steps for each dataset are outlined in the following sections.

A. LISTINGS DATA PRE-PROCESSING

The Listings dataset includes an amenities column that lists key amenities for each property. To make it easier for users to choose accommodations, these amenities were divided into distinct categories. There are 77 unique amenities, which were grouped into 9 broader categories. For example, amenities like 'Hot shower,' 'Shower gel,' 'Hair dryer,' 'Shampoo,' 'Conditioner,' and 'Body Soap' were categorized under 'Bath Products,' while 'Oven,' 'Hot water kettle,' 'Cooking basics,' and 'Microwave' were categorized under 'Kitchen Appliances.' If a listing lacked any amenities from the 9 categories or did not have host verification information, the missing values were filled with zeros. Figure 1 illustrates the percentage of NaN values in each column of the Listings dataset.

Next, it was crucial to evaluate the number of NaN values in each feature column. Figure 1 shows the percentage of NaN values in each column of the Listings dataset. The analysis identified that the columns 'neighbourhood group cleansed,' 'license,' 'bathrooms,' and 'calendar updated' contained NaN values, leading to their exclusion from further analysis. Additionally, the seven Review Scores columns, which are our

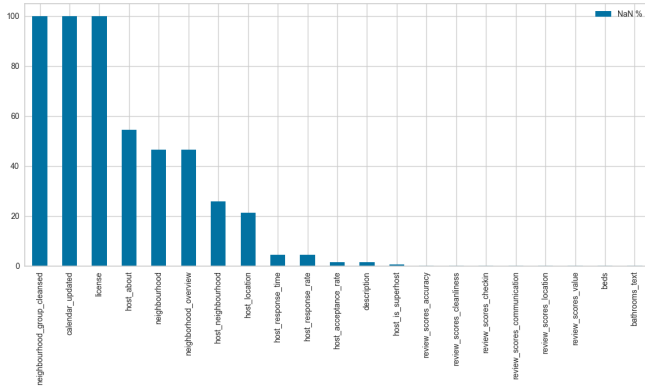


Fig. 1. %NaN values present in each feature columns from Listings dataset

target variables, had about 20% NaN values. Since these are essential for our model, imputing these values could introduce bias, so any rows with NaN values were removed from the dataset.

A detailed examination of the ‘price’ feature revealed that prices were listed in dollars (\$). To ensure machine-readability, I removed the ‘\$’ and ‘,’ symbols and converted the string values to integers. This preprocessing step was vital for maintaining data consistency. The ‘price’ feature displayed numerous outliers, with an average price of about \$167 and a median of \$110, while the maximum price reached \$1,600. As shown in Figure 2, most prices fell within the \$50 to \$300 range, leading to the removal of rows with prices outside this range.



Fig. 2. Plot of Price Feature

Finally, the distribution of host response rates was plotted to understand the overall picture of host responsiveness in the dataset. As shown in Figure 3, the majority of hosts have a highly concentrated response rate of 100%, indicating they respond to guest inquiries very promptly. A small number of hosts have response rates between 0% and 20%, showing significantly lower responsiveness. There is also a smaller group of hosts with moderate response rates between 20% and 80%. This right-skewed distribution may result from the platform’s rules or incentives, hosts’ sense of responsibility, and potential biases in data collection. Overall, this high response rate enhances the platform’s user experience and satisfaction, but improvement measures are needed for the few hosts with low response rates.

The rating values for the seven target variables were closely

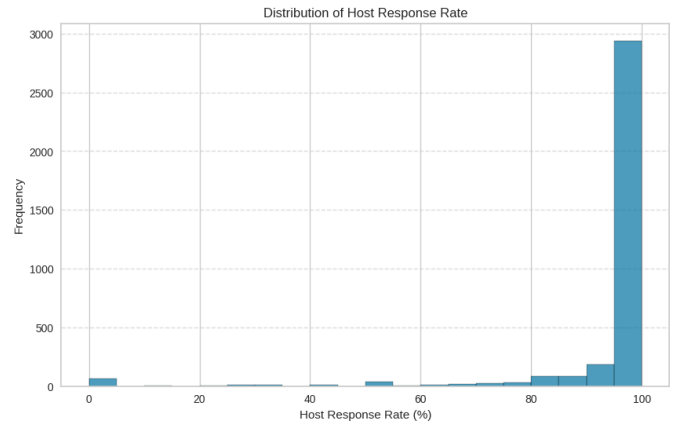


Fig. 3. Plot of Host Response Rate

clustered. Figure 4 illustrates the distribution of all rating columns, which predominantly range between 4 and 5. To address this clustering, I used a Binning Approach, categorizing continuous data into discrete bins based on quantile ranges. This method helps avoid bias by ensuring equal-sized bins. As a result, we predict the rating bin for a given listing rather than the exact numeric rating, transforming this into a classification problem where the model predicts a categorical outcome based on input parameters.

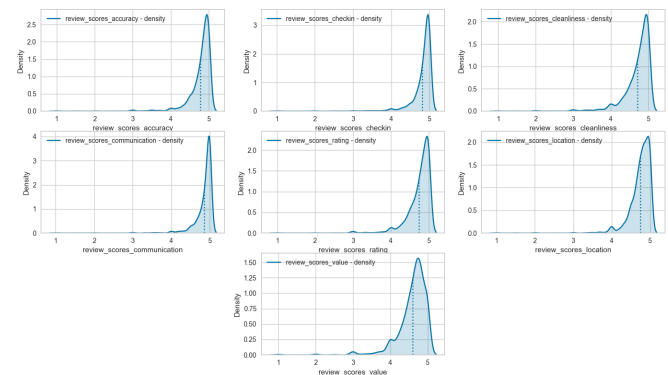


Fig. 4. Distribution Plot of all the Target Columns (Review Rating Columns)

In the Binning Approach, continuous data is categorized into discrete bins based on specific criteria, which facilitates further analysis. For this use case, I chose to bin individual rating values based on their quantile range. The primary reason for using this approach is that the quantile method allows for the creation of equal-sized bins.

By categorizing a given rating value into equal-sized bins according to its quantile range, we avoid bias where one bin has more records than others, making this method suitable for our problem. Once the values are binned, we predict the rating bin for a given listing instead of the exact numeric rating, transforming this into a classification problem where the model predicts a categorical variable based on input parameters.

To categorize ratings into discrete bins, I used the Pandas

DataFrame QCut functionality, which segments the data into equal-sized bins based on distribution percentiles. Ratings were first converted to percentages before binning. For features like ‘host response rate’ and ‘host acceptance rate,’ I removed the ‘%’ signs and converted the resulting figures into numeric format for further processing. Features such as ‘host response time,’ ‘host is superhost,’ ‘host identity verified,’ ‘instant bookable,’ ‘room type,’ ‘neighbourhood cleansed,’ and ‘has availability’ were label-encoded using sklearn’s preprocessing library to ensure they are machine-readable. For scaling, I applied min-max scaling, which preserves the original data distribution. In min-max scaling, each data point is scaled to a new value using the formula.:

$$x_{scaled} = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

B. REVIEWS DATA PRE-PROCESSING

The Reviews dataset includes features such as listing ID, ID, date, reviewer ID, reviewer name, and comments. With 230,065 reviews, a single listing may have multiple reviews. For analysis, I focused solely on the listing ID and comments. The comments feature contained non-English entries and some comments composed entirely of special characters, emojis, or HTML tags. To make these features usable, I needed to preprocess them by removing emojis and special characters, while retaining only English reviews.

Initially, I eliminated all review rows lacking associated comments. Then, I filtered for English reviews using the **fasttext** Python library and its pretrained model, which identifies the language of any text. Next, I removed emojis from the text using the **emoji** Python library. Each comment was analyzed character by character to check for matches with the emoji data. If a comment consisted entirely of emojis or was predominantly composed of them, it was discarded. I set a threshold of 20 characters for comments to ensure they were substantial enough for model training. After preprocessing, 220,551 reviews were retained out of the original 230,065, resulting in a 4% reduction of irrelevant reviews. The total number of comments for a given listing was updated after preprocessing was completed.

1) **TF-IDF METHODOLOGY:** To utilize the review text effectively, preprocessing was necessary to make it machine-readable. I implemented the Term Frequency-Inverse Document Frequency (TF-IDF) methodology, which assesses the significance of a word within a text collection. The TF-IDF of a term, such as a review for a listing, by considering both the frequency of the term in the document and its rarity across all documents.

TF-IDF of term t in a document d is defined as:

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{idf}(t)$$

where:

$$\text{tf}(t, d) = \frac{\text{terms } t \text{ occur in doc } d}{\text{total terms in } d}$$

$$\text{idf}(t) = 1 + \log \left(\frac{1 + \text{documents in corpus}}{1 + \text{df}(t)} \right)$$

Here, $\text{tf}(t, d)$ represents the term frequency, and $\text{idf}(t)$ represents the inverse document frequency.

The rationale for using TF-IDF in this context was to extract important tokens from each review to serve as features in the machine learning model. After the preprocessing steps outlined in Section 2.1, I removed all stopwords from the review comments. For the TF-IDF implementation, I utilized the ‘TfidfVectorizer’ from sklearn’s text feature extraction library, limiting the maximum number of features returned to 33. This limitation helps mitigate potential bias and overfitting, as the distribution of word counts per comment varies significantly.

The limitation on the number of features returned by ‘TfidfVectorizer’ serves two main purposes. First, considering the wide range of word counts per comment, there are significant outliers—from very short comments of 10-15 words to lengthy ones spanning 500-1000 words. Based on the distribution shown in Figure 5, most data falls within the 20-40 words per comment range. The statistical summary of the word counts per comment shows a mean of approximately 44 and a median of 33. Restricting the maximum features to the median aligns with the central tendency of the data distribution.

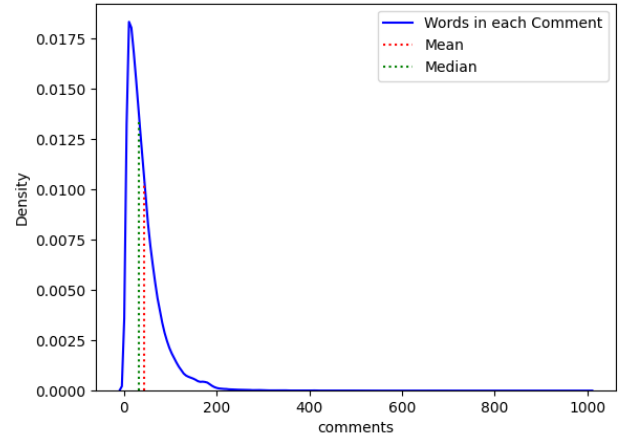


Fig. 5. Plot of Words in each comment

Moreover, an increase in the number of features could potentially bias the data and lead to overfitting, resulting in high training accuracy but poor testing accuracy.

Finally, after computing TF-IDF for all review comments, I calculated the mean of all features created by ‘TfidfVectorizer’, grouped by Listing ID. This approach ensures that each listing has only one row, facilitating seamless integration with the Listings dataset.

C. FEATURE SELECTION AND IMPORTANCE

After cleaning and processing the Listings and Reviews datasets, I aimed to identify features that significantly impact the machine learning model. Using ‘review scores accuracy’ as an example, I calculated the correlation coefficients between the features and the target variables. A threshold of 0.01

was set, retaining only features with correlation coefficients exceeding this value.

As shown in Figure 6, features such as 'host is superhost' (correlation coefficient of 0.257), 'host listings count' (0.232), and 'host total listings count' (0.215) have a strong positive correlation with review scores accuracy. Additionally, positive terms used in reviews, such as 'good' (0.211), 'home' (0.166), and 'recommend' (0.166), also show high correlations, indicating that the content of the review influences the rating.

Regarding amenities and services, features such as parking (0.097), food service (0.078), and security (0.079) displayed notable correlations. Host response rate (0.070) and acceptance rate (0.033) also ranked high, showing the impact of hosts' timely responses and acceptance of booking requests on ratings.

This analytical approach helped identify multiple features significantly affecting ratings, providing key inputs for subsequent machine learning models and enhancing the accuracy and reliability of predictions.

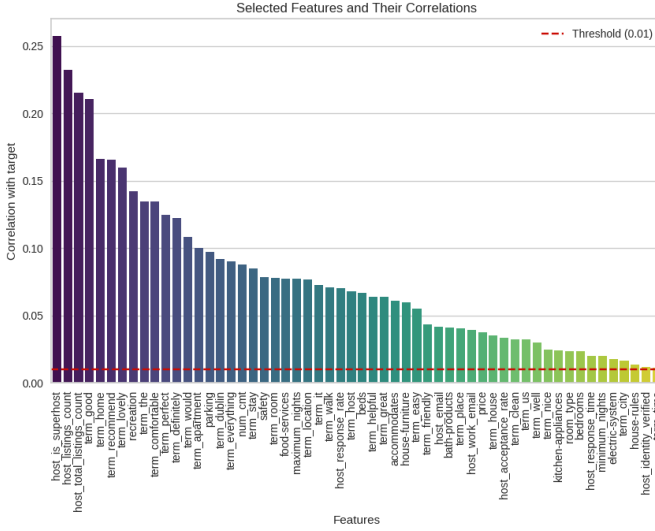


Fig. 6. The Correlations of Features - eg.review_scores_accuracy

Figure 6 shows that for a listing to achieve a very high rating (Bin 2), the most significant terms in the review comments are 'Home,' 'Perfect,' and 'Recommend.' For a rating in Bin 1, reviews must include terms like 'Would,' 'Lovely,' or 'Home,' along with amenities like Kitchen Appliances. The analysis also highlights differences in the intensity of the term 'Home' and the number of comments per bin, especially when comparing Superhosts to non-Superhosts. This suggests that these features are crucial in developing the model.

Figure 7 shows that the intensity of the term 'Home' and the number of comments per bin differ across all three review bins. The number of records in each review bin varies significantly when the host is a Superhost compared to when the host is not a Superhost. This highlights the importance of these features in model development. A similar analysis was conducted for all the features, and the final model was built using a total of 60 features.

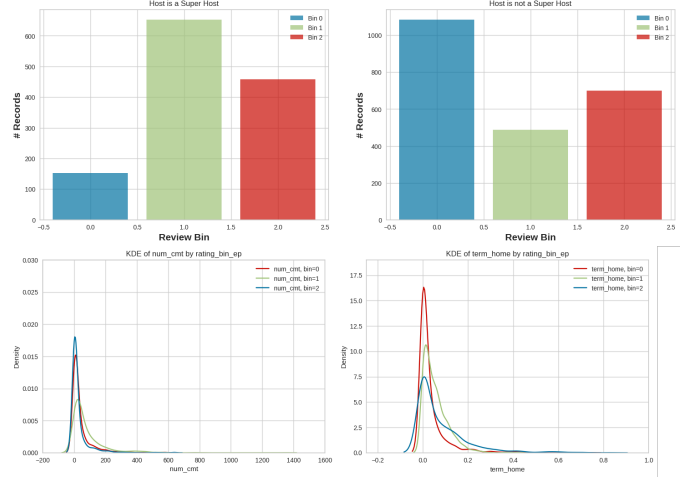


Fig. 7. Analysis Important Features

III. MODELS OF MACHINE LEARNING

Given that this use case is framed as a classification problem, the initial models to test and implement include Logistic Regression, k-Nearest Neighbors (k-NN), and Decision Tree Classifier. All models were executed as part of this assignment, and their results were compared to determine the best fit. The combined dataset of Listings and Reviews was split into an 80-20 ratio, with 3741 data rows. The models were trained on 78% of the dataset and tested on the remaining 20%.

A. LOGISTIC REGRESSION AND PARAMETER TUNING

Logistic Regression is a supervised learning algorithm that uses labeled data for training. During this process, the model assigns weights (coefficients) to each feature, optimizing these weights to minimize error through a cost function. To prevent overfitting, a penalty (L1, L2, or None) is applied as a hyperparameter. For this case, I explored a range of C values, controlling the regularization strength from 0.001 to 1000. The C values tested were [0.001, 0.1, 1, 10, 100, 1000], exposing the model to varying levels of regularization, where smaller C values indicate stronger regularization.

Given the multiclass nature of the problem, where the model predicts which review bin an Airbnb listing falls into (0, 1, or 2), I selected the 'newton-cg' solver recommended by scikit-learn for multinomial loss problems. This solver supports the L2 penalty and computes the Hessian matrix, making it suitable for this scenario. To identify the optimal C value for each target variable, I plotted the cross-validation (CV = 5) graph. Figure 8 presents accuracy scores across varying C values for the review scores target variable. The results indicate that increasing C (reducing regularization) generally enhances model accuracy, although improvements diminish after a certain point (C = 10). The optimal C value and corresponding accuracies for each target variable are summarized in Table I for degree 1 and Table II for degree 2.

For most target variables, the optimal value of C is 10 or less for degree 1 models and similar for degree 2 models.

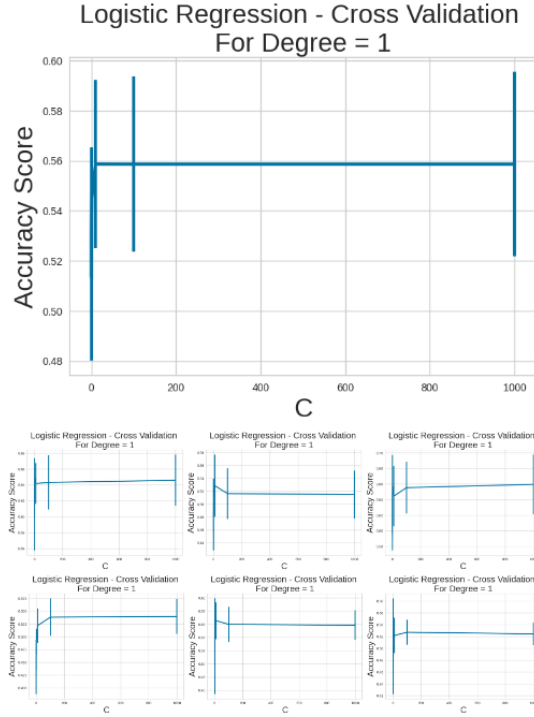


Fig. 8. Cross-Validation Accuracy Scores for Degree 1

TABLE I
VALUES OF C FOR DIFFERENT TARGET VARIABLES (DEGREE=1)

Target Variable	C Value	Train Accuracy	Test accuracy
accuracy	10	0.60	0.59
checkin	10	0.69	0.68
cleanliness	100	0.58	0.54
communication	10	0.71	0.68
location	10	0.60	0.60
rating	100	0.60	0.58
value	100	0.60	0.60

Increasing the value of C further reduces the regularization penalty, risking overfitting. This is reflected in the optimal C values shown in Tables I and II. Transforming the training dataset into a polynomial feature space of degree 2 was computationally expensive and did not significantly improve results, as there was still a notable gap between the training and testing accuracy scores.

When examining the results more closely, it is apparent that for certain variables like check-in and communication, higher C values (such as 10) yielded better generalization performance, indicating that moderate regularization was beneficial. Conversely, for variables like cleanliness and value, very high C values (such as 100 or 1000) were optimal, suggesting these features were more sensitive to regularization and benefitted from stronger penalization to avoid overfitting.

Additionally, I augmented and transformed the training dataset into a polynomial feature space of degree 2 to predict the review bins a given listing might belong to. This process yielded slightly better results in some cases but was generally not significantly better than the degree 1 models, as indicated

TABLE II
VALUES OF C FOR DIFFERENT TARGET VARIABLES (DEGREE=2)

Target Variable	C Value	Train Accuracy	Test accuracy
accuracy	10	0.70	0.60
checkin	10	0.77	0.68
cleanliness	100	0.64	0.55
communication	10	0.71	0.68
location	10	0.70	0.60
rating	100	0.65	0.60
value	1000	0.65	0.58

by the results in Table II. The increase in computational complexity did not justify the marginal improvement in accuracy. This was evident in variables like check-in and location, where both degree 1 and degree 2 models performed similarly.

The cross-validation plot for all the target variables predicted by the model is presented in Figure 9, illustrating how the model's performance varied across different regularization strengths and polynomial degrees. The plots show that the model's accuracy tends to stabilize or even decline beyond a certain C value, reinforcing the importance of selecting an appropriate regularization parameter.

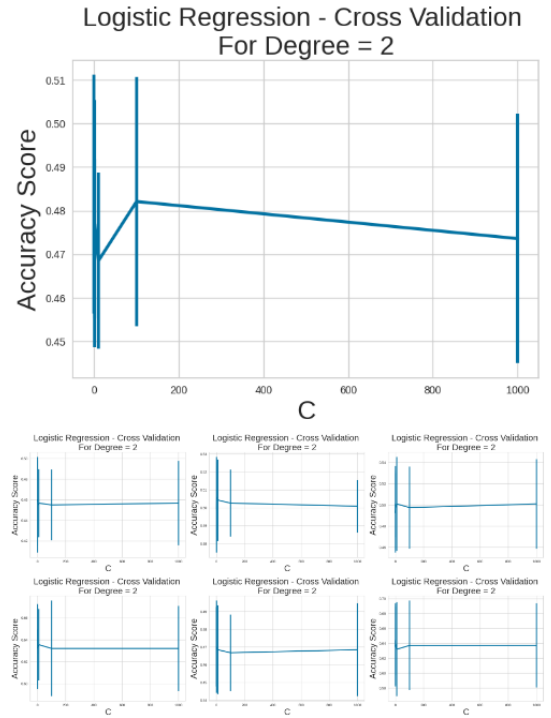


Fig. 9. Cross-Validation Accuracy Scores for Degree 2

B. k -NN CLASSIFIER AND PARAMETER TUNING

The k -Nearest Neighbors (k -NN) algorithm is a supervised learning method that classifies data based on their target variables or classes. The model uses its hyperparameter (k neighbors) and the distance between classes to classify new data points. In this analysis, I varied the value of k while keeping weights uniform, selecting only odd numbers to avoid

ties. I implemented 5-fold cross-validation to assess model performance, testing k values in the range of [1, 3, 5, 7, 9, 11, 13, 15, 17, 19].

The performance of the k-NN classifier depends significantly on the choice of k, the number of neighbors considered. A smaller k can lead to a model that is too complex, capturing noise in the dataset (overfitting), while a larger k can oversimplify the model, missing important trends (underfitting). To find the optimal k value, I conducted experiments across a range of odd values to avoid ties during classification.

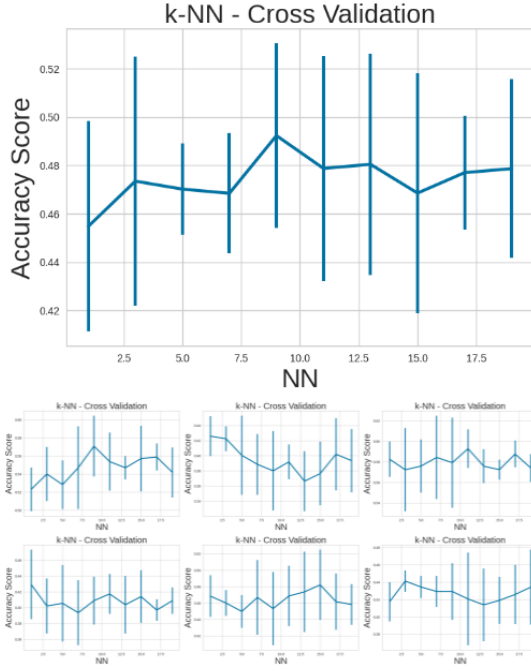


Fig. 10. k-NN Cross Validation for the target variable

Figure 10 illustrates the cross-validation plot for review scores accuracy across different k values. The analysis showed that the model reached its highest accuracy with k set to 7. Although k=1 produced similar accuracy, it was found to be too complex and hard to interpret. Increasing k simplifies the model, reducing dataset noise. Therefore, the optimal k value for review scores accuracy is 7, with training accuracy at 63.6% and testing accuracy at 49.6%. Similarly, the optimal k values for other target variables demonstrated a trend of balancing model complexity and performance.

The detailed results for each of the seven target variables (accuracy, checkin, cleanliness, communication, location, rating, and value) are summarized in Table III. These results highlight the importance of selecting an appropriate k value for each specific target to achieve optimal performance.

C. DECISION TREE AND PARAMETER TUNING

The Decision Tree algorithm is a supervised learning method that divides data based on their features to make decisions. This analysis examines the performance of decision tree classifiers with varying tree depths (from 3 to 19) across

TABLE III
VALUE OF HYPER-PARAMETER K WITH THE CORRESPONDING TRAIN AND TEST ACCURACY FOR ALL THE 7 TARGETS

Target Variable	K Value	Train Accuracy	Test accuracy
accuracy	7	0.63	0.50
checkin	19	0.67	0.61
cleanliness	7	0.60	0.44
communication	5	0.75	0.62
location	7	0.60	0.43
rating	7	0.64	0.50
value	19	0.53	0.46

multiple target variables: accuracy, check-in, cleanliness, communication, location, rating, and value. Each tree depth was evaluated using 5-fold cross-validation to assess model performance, testing depths in the range of [3, 5, 7, 9, 11, 13, 15, 17, 19].

The performance of the Decision Tree classifier is significantly influenced by the depth of the tree. A shallow tree (small depth) might underfit the data, missing important patterns, while a deep tree (large depth) might overfit, capturing noise in the dataset. To find the optimal tree depth, experiments were conducted across a range of depths.

From the Figure 11, for the accuracy target, a depth of 3 provided a train accuracy of 61% and a test accuracy of 58%, while a deeper tree at depth 7 resulted in a train accuracy of 74% but a lower test accuracy of 56%, indicating overfitting.

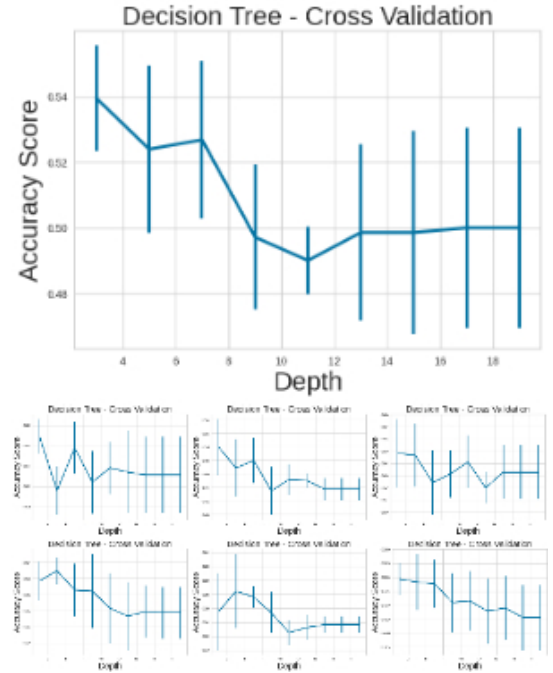


Fig. 11. Decision Tree Cross Validation for the target variable

The detailed results for each of the seven target variables highlight the importance of selecting an appropriate tree depth to achieve optimal performance. Table IV summarizes the value of hyper-parameter depth with the corresponding train and test accuracy for all the 7 targets:

TABLE IV
VALUE OF HYPER-PARAMETER DEPTH WITH THE CORRESPONDING TRAIN
AND TEST ACCURACY FOR ALL THE 7 TARGETS

Target Variable	Depth Value	Train Accuracy	Test accuracy
accuracy	3	0.61	0.58
checkin	5	0.74	0.70
cleanliness	5	0.64	0.52
communication	5	0.70	0.67
location	5	0.59	0.54
rating	5	0.69	0.60
value	7	0.71	0.55

The figures and tables illustrate the impact of tree depth on model performance, emphasizing the importance of choosing the right depth to balance model complexity and accuracy.

D. RESULT EVALUATION AND MODELS COMPARISON

In this section, the performance of the Logistic Regression, k-NN, and Decision Tree models will be evaluated and compared to determine which model best fits the dataset and achieves the desired accuracy in predicting Airbnb listing ratings. The ROC-AUC (Receiver Operating Characteristic - Area Under the Curve) metric is used to assess the models. An ROC curve is a plot of True Positive Rate vs. False Positive Rate, while AUC is the area under the ROC curve, providing a single value for each class rather than a curve. For an ideal classifier, the ROC curve gives a point in the top left corner, indicating 100% True Positives and 0% False Positives; similarly, for an ideal classifier, $AUC = 1$, and a random classifier has $AUC = 0.5$.

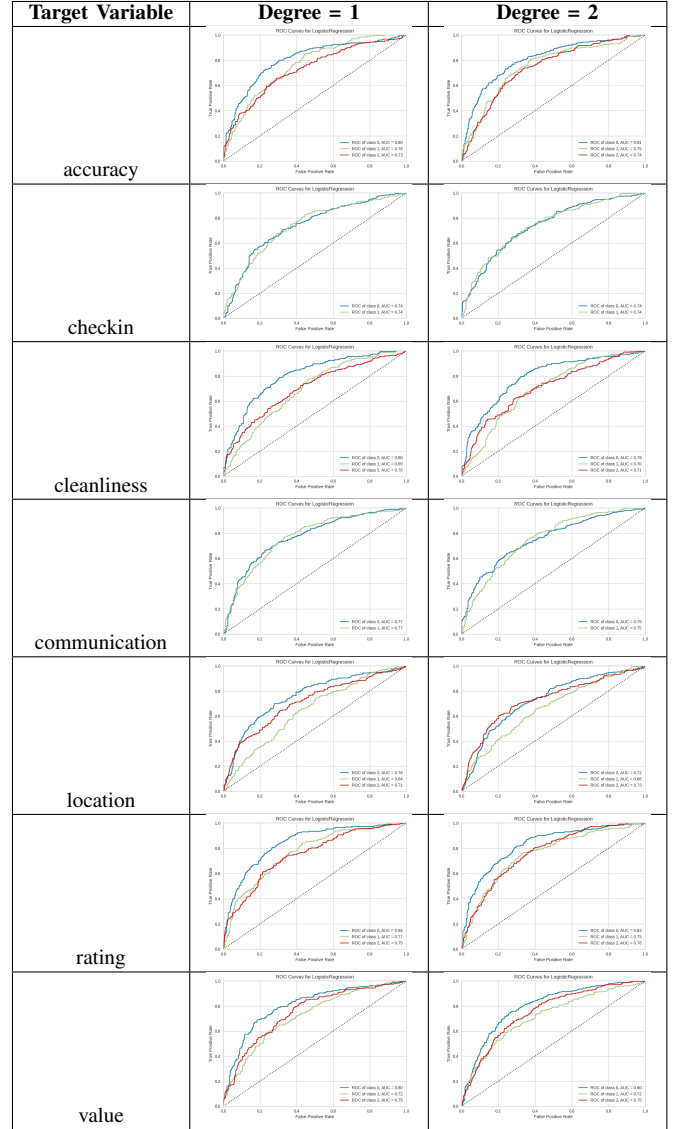
Please find the ROC-AUC curves of Logistic Regression, k-NN, and Decision Tree models for various target variables in the following Table VI and Table ??.

From the values presented in the ROC-AUC curves, it is evident that the Logistic Regression model consistently outperforms both the k-NN and Decision Tree classifiers for all target variables. Comparing Logistic Regression with degree=1 and degree=2, we observe that the performance is relatively similar, with degree=1 slightly outperforming degree=2 in some instances, and vice versa in others. For instance, in the 'review_scores_accuracy' target variable, degree=1 achieves an AUC of 0.80, 0.76, and 0.73 for classes 0, 1, and 2 respectively, while degree=2 achieves 0.81, 0.75, and 0.74. Both degrees of Logistic Regression are robust and significantly better than k-NN and Decision Tree models.

The AUC for all classes across all target variables is higher for Logistic Regression compared to k-NN and Decision Tree classifiers. For example, in the 'review_scores_cleanliness' target variable, Logistic Regression (degree=1) achieves an AUC of 0.80, 0.69, and 0.70 for classes 0, 1, and 2 respectively, whereas k-NN achieves 0.62, 0.62, and 0.63, and Decision Tree achieves 0.65, 0.60, and 0.61.

The Dummy Classifiers further emphasize the superiority of the selected models. The Frequent Class Dummy achieved an accuracy of 36.2%, with all predictions skewed towards the most frequent class, and the Uniform Class Dummy achieved an accuracy of 33.8%, with predictions uniformly distributed

TABLE V
ROC-AUC FOR DIFFERENT DEGREE OF LOGISTIC REGRESSION



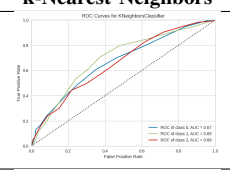
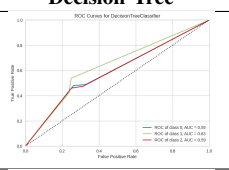
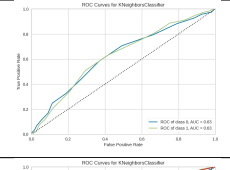
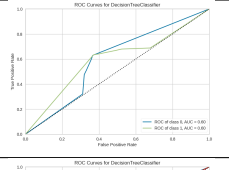
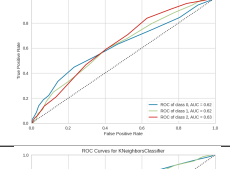
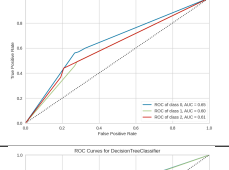
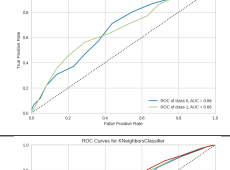
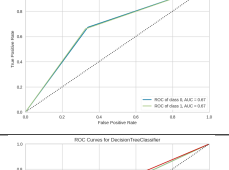
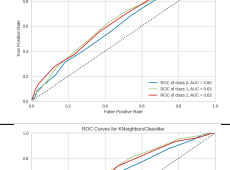
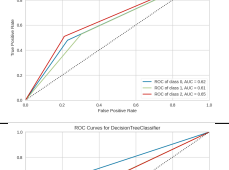
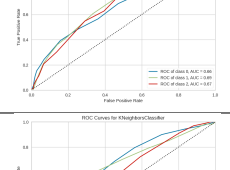
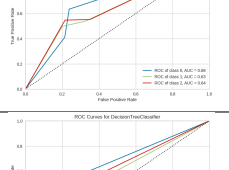
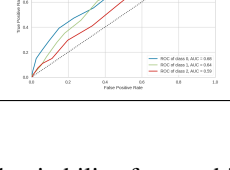
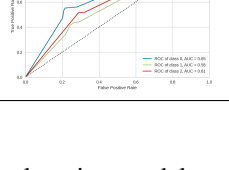
across classes. Both of these baselines perform significantly worse than the Logistic Regression, k-NN, and Decision Tree models.

Therefore, Logistic Regression, regardless of the degree, is the most suitable model for this dataset and problem statement. The selected models' performance is optimal for this use case, as evidenced by their superior AUC values and accuracy compared to the baseline classifiers.

IV. CONCLUSION

In this research, our goal was to forecast possible ratings for Airbnb listings in Dublin by examining different features and services offered. We thoroughly processed the data and engineered features, categorizing amenities, addressing missing values, and effectively utilizing review text data. We carefully prepared datasets from the Inside Airbnb website to maintain

TABLE VI
ROC-AUC BETWEEN K-NEAREST NEIGHBORS AND DECISION TREE

Target Variable	k-Nearest Neighbors	Decision Tree
accuracy		
checkin		
cleanliness		
communication		
location		
rating		
value		

data integrity and suitability for machine learning models.

We utilized three main machine learning models in our analysis: Logistic Regression, k-Nearest Neighbors (k-NN), and Decision Tree Classifiers. Each model underwent thorough testing and tuning to identify the optimal hyperparameters for the best predictive performance. The results from these models revealed several important insights:

- 1) **Logistic Regression:** Logistic Regression proved to be the strongest model for all target variables, consistently achieving higher accuracy and ROC-AUC scores compared to k-NN and Decision Tree models. The best regularization parameter (C) differed across target variables, with moderate regularization typically producing the best outcomes. Degree 1 polynomial features performed slightly better than degree 2, suggesting that adding complexity did not notably improve predictive accuracy.

2) **k-Nearest Neighbors (k-NN):** The k-NN model highlighted the significance of choosing the right k value. Striking a balance between underfitting and overfitting was essential, with k=7 yielding the best results for most target variables. Although k-NN is simple and easy to implement, it was less effective than Logistic Regression, especially for higher-dimensional data.

3) **Decision Tree Classifier:** The performance of the Decision Tree model was highly dependent on tree depth. Shallow trees tended to underfit, while deep trees often overfitted the data. The optimal tree depths varied for different target variables, but overall, the Decision Tree model did not perform as well as Logistic Regression in terms of predictive accuracy.

This research offers substantial academic and practical benefits. Academically, it showcases the application of various machine learning algorithms to real-world problems, particularly in predicting Airbnb listing ratings. The detailed feature engineering and model evaluation processes serve as a valuable guide for future studies aiming to predict categorical outcomes using complex datasets.

Practically, the study provides actionable insights for Airbnb hosts and platform developers. By pinpointing key features that affect listing ratings, hosts can focus on enhancing specific aspects of their services, such as cleanliness, communication, and amenities, to boost guest satisfaction. Additionally, the platform can use these predictive models to offer personalized recommendations and improve the overall user experience.

V. ANSWER OF ASSIGNMENT 2

A. QUESTION 1

Following are the two situations where Logistic Regression would give inaccurate results:

1) Data Imbalance

If the dataset consists of imbalanced data for the target variable, as shown in Figure 12 left. The target value in this data is made up of 34 "1s" and just 6 "0s." As seen in Figure 12 right, the model can only predict for the majority class, or 1, when predicting this feature.

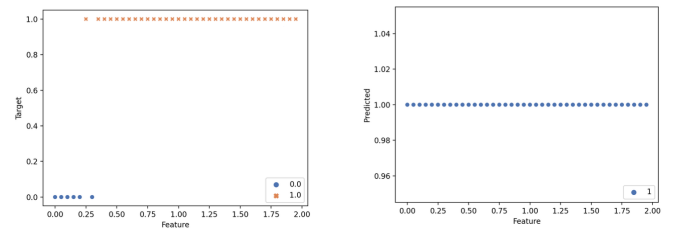


Fig. 12. Plot of Data Imbalance

2) No Linear Correlation

The model cannot correctly forecast the target variable if there is no linear correlation between the data. For instance, I arbitrarily generated two features and one target variable. The two features and their corresponding

goal values are shown in Figure 13 left. The following plot, Figure 13 right, shows that the values are predicted linearly rather than scatteredly as they were formed when estimating the target values from the attributes.

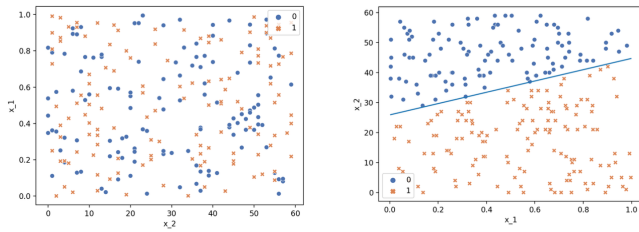


Fig. 13. Two Features and Target Values

B. QUESTION II

1) KNN's Advantages

To identify the data points that are part of a particular class, the model uses a simple calculation. The K data points that are closest to the new data point are chosen by the algorithm after it has calculated the distance between the new data point and every other data point in the dataset. Moreover, the prediction process is dependent on the complete dataset. Without more training, the method can be applied simply to a new dataset. Finally, it can be applied to regression as well as classification. The KNN model uses a majority vote from its neighbors to classify the target variables. The property value of an item is the result of regression. The average of its K Nearest Neighbors is this value.

2) KNN's Disadvantages

Because it scans the complete dataset to find the N-Neighbors, it may take a while to forecast the target class for a new data instance. Additionally, the technique is susceptible to anomalies or irrelevant data in the dataset. To deal with these anomalies, the data must be scaled properly. Finally, a prediction is the only step in the creation of a KNN algorithm model. When training data is vast and prediction time is crucial, there is a significant overhead.

3) MLP's Advantages

MLP classifiers are capable of modeling intricate relationships found in data. Non-linear correlations between the input and output variables can be learned by these classifiers. They may also process high-dimensional data, such as text, audio, and photos, by teaching the hidden layers to learn distributed representations of the input data. In addition, MLP classifiers may be trained on big datasets and excel in a variety of real-world classification tasks with respect to predicted accuracy. Given that the model is stored once it has been trained on data, they can be utilized to rapidly forecast the target values. The learning rate, activation function, number of hidden layers, degree of regularization, dropout, and

batch size are among the hyperparameters in MLP that can be adjusted.

4) MLP's Disadvantages

Training a large dataset can be time-consuming, particularly if the model has a lot of hidden layers. The model has numerous hidden layers, therefore fine-tuning its hyperparameters is necessary to achieve acceptable results. Additionally, as the number of hidden layers in an MLP model increases, the model gets increasingly complex and difficult to understand. As a result, using it as a black box model obscures the details of how it functions.

C. QUESTION III

The dataset is resampled several times in K-Fold Cross Validation in order to train and assess the model on various subsets of the data and determine how well the model generalizes. A model's performance may be overestimated if it is trained on a single dataset split and then performs poorly on additional, unknown data. A more accurate estimation of the model's performance on unseen data with various data splits can be obtained by resampling the data several times across different folds. For instance, I used np.arange to generate a feature X with values ranging from 0 to 1, and I used feature X to build a target Y. The training data is split five times at random, and each time the desired output is predicted, Figure 14 is produced. For each of these separate splits, or folds, the Mean Squared Error is determined. The MSE varies for every fold, as can be observed. It is acknowledged that using a single train-test split to evaluate the model is not a recommended approach.

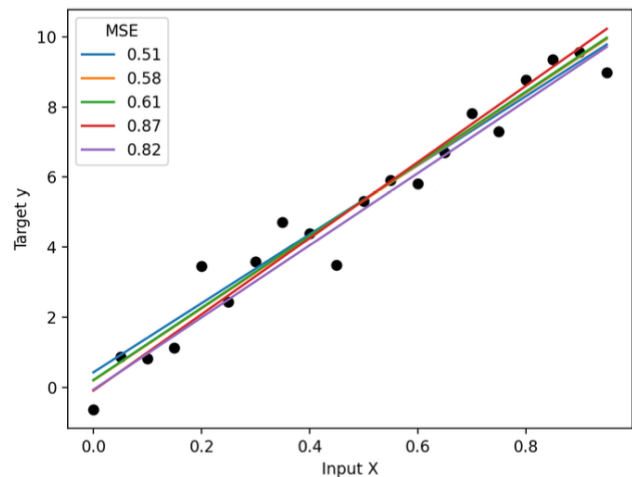


Fig. 14. Random data splits MSE

Selecting the Value of K: The model would be trained and assessed on a restricted number of folds if a small value of k was chosen. This would result in a significant bias, which would suggest that the model makes assumptions about the target variable. However, choosing a large number for k will cause the model to be trained numerous times, which will

result in a high variance. This means that the model will overfit the data and become sensitive to even slight changes in the target variables.

D. QUESTION VI

Q is the future value that we wish to anticipate for a time series data set at time $t+q$. Lagged output values can be utilized to generate features for the time series at t , $t-1$, and $t-2$. These values provide useful information for making future value predictions. For instance, let's say we wish to develop a times series feature that allows us to forecast how many products would be used each month. The current value is shifted by one or two steps to produce the lag values. The data for feature tag_1 is moved by one step, as shown in Figure 15, and the data for tag_2 is shifted by two steps.

	Date	Product Usage	tag_1	tag_2
0	2024-07-15	1000	NaN	NaN
1	2024-08-15	2000	1000.0	NaN
2	2024-09-15	4000	2000.0	1000.0
3	2024-10-15	8000	4000.0	2000.0
4	2024-11-15	16000	8000.0	4000.0
5	2024-12-15	32000	16000.0	8000.0

Fig. 15. Time Series Features for Product Usage

APPENDIX

```

1 import warnings
2 warnings.filterwarnings('ignore')
3 import pandas as pd
4 from sklearn.feature_extraction.text import
   TfidfVectorizer
5 from nltk.corpus import stopwords
6 import pandas as pd
7 import numpy as np
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 import regex
11 import emoji
12 from pycountry import languages
13 import fasttext
14
15 def split_count(info):
16
17     emoji_list = []
18     data = regex.findall(r'\X', info)
19     for word in data:
20         if any(char in emoji.EMOJI_DATA for char in
21            word):
22             emoji_list.append(word)
23
24     return len(emoji_list)
25
26 reviews_raw = pd.read_csv("reviews.csv")
27 PRETRAINED_MODEL_PATH = 'lid.176.bin'
28 model = fasttext.load_model(PRETRAINED_MODEL_PATH)
29 reviews_raw['lang'] = range(0, len(reviews_raw))
30
31 for index, row in reviews_raw.iterrows():
32     if type(row['comments']) == type('s'):
33         predictions = model.predict(row['comments'])
34         l = predictions[0][0].split('_label_')[1]
35         if l != 'ceb' and l != 'nds' and l != 'war'
           and l != 'wuu':

```

```

36         reviews_raw['lang'][index] = 1
37
38 reviews_raw = reviews_raw[reviews_raw['comments'].
   notna()]
39
40 reviews_raw['emoji_count'] = reviews_raw.comments.
   apply(split_count)
41
42 reviews_raw['lenstr'] = reviews_raw['comments'].str.
   len()
43
44 reviews_raw['test_emoji'] = (((reviews_raw['
   emoji_count'] == reviews_raw['lenstr']) |
45                               (2*reviews_raw['
   emoji_count'] == reviews_raw['lenstr'])) & (
46                               reviews_raw['emoji_count'] != 0))
47
48 reviews_raw = reviews_raw[reviews_raw.test_emoji ==
   False]
49
50 reviews_raw = reviews_raw.drop(columns='test_emoji')
51 reviews_raw = reviews_raw.drop(columns='lenstr')
52 reviews_raw = reviews_raw.drop(columns='emoji_count'
   )
53
54 reviews_raw.comments = reviews_raw.comments.apply(
   lambda x: x.encode('ascii', 'ignore').decode('
   ascii'))
55
56 reviews_raw['emoji_count'] = reviews_raw.comments.
   apply(split_count)
57
58 reviews_raw = reviews_raw[(reviews_raw.comments.str.
   len() > 20)]
59
60 reviews_raw['comments'] = reviews_raw['comments'].
   str.replace('<br/>', '')
61
62 reviews_raw = reviews_raw[reviews_raw.lang == '_en']
63
64 reviews_raw.to_csv("reviews_processed.csv")
65
66 reviews = pd.read_csv("reviews_processed.csv",
   lineterminator='\n')
67
68 print("Mean: ", reviews['comments'].str.split().str.
   len().mean())
69
70 print("Median: ", reviews['comments'].str.split().
   str.len().median())
71
72 i = sns.kdeplot(reviews['comments'].str.split().str.
   len(), color="b", label='Words in each Comment')
73
74 xi = i.lines[0].get_xdata()
75 yi = i.lines[0].get_ydata()
76
77 meani = reviews['comments'].str.split().str.len().
   mean()
78
79 mediani = reviews['comments'].str.split().str.len().
   median()
80
81 heighti = np.interp(meani, xi, yi)
82
83 heighti2 = np.interp(mediani, xi, yi)
84
85 i.vlines(meani, 0, heighti, color='r', ls=':', label
   ='Mean')
86
87 i.vlines(mediani, 0, heighti2, color='g', ls=':',
   label='Median')
88
89 plt.legend()
90
91 plt.show()
92
93 stop = stopwords.words('english')
94
95 reviews['comments_stp_rem'] = reviews['comments'].
   apply(lambda x: ' '.join([word.lower() for word
96    in x.split()
97
98    if word not in (stop)]))
99
100 tr_idf_model = TfidfVectorizer(analyzer = 'word',
   max_features=33)

```

```

85 tf_idf_vector = tr_idf_model.fit_transform(reviews.
      comments_stp_rem)
86
87 tf_idf_array = tf_idf_vector.toarray()
88 words_set = tr_idf_model.get_feature_names_out()
89 df_tf_idf = pd.DataFrame(tf_idf_array, columns =
      words_set)
90 df_tf_idf['listing_id'] = reviews.listing_id
91
92 test = df_tf_idf.groupby('listing_id').size()
93 review_1 = pd.DataFrame({'listing_id':test.index, '
      num_cmt':test.values})
94
95 test2 = pd.DataFrame([])
96 for i in words_set:
97     test2[f'term_{i}'] = df_tf_idf.groupby(['
      listing_id']) [f'{i}'].mean()
98
99 review_final = review_1.merge(test2, on='listing_id'
      )
100
101 review_final.to_csv("reviews_final.csv")

```

```

1 import warnings
2 warnings.filterwarnings('ignore')
3
4 import pandas as pd
5 import numpy as np
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 from sklearn import preprocessing
9 from sklearn.linear_model import Lasso,
      LogisticRegression, Ridge
10 from sklearn.model_selection import train_test_split
11 from sklearn.metrics import mean_squared_error
12 from sklearn.preprocessing import MinMaxScaler
13 from sklearn.metrics import mean_squared_error,
      accuracy_score, log_loss
14 from sklearn.dummy import DummyRegressor,
      DummyClassifier
15 from yellowbrick.classifier import ROCAUC
16 from sklearn.metrics import confusion_matrix,
      precision_score, recall_score, auc
17 from sklearn.metrics import f1_score,
      classification_report, roc_curve
18 from sklearn.neighbors import KNeighborsClassifier
19 from sklearn.tree import DecisionTreeClassifier
20 from sklearn.preprocessing import PolynomialFeatures
21 from sklearn import metrics
22 from sklearn.model_selection import cross_val_score
23
24 def processing_and_merge(listings, reviews):
25
26     amenities = listings['amenities'].str.split(',',
      expand=True)
27     amenities = amenities.loc[amenities[77].notnull()]
28     amenities_list = amenities.iloc[0]
29     amenities_list = amenities_list.to_list()
30
31     listings.loc[listings['amenities'].str.contains('
      Hot Water|Shower gel|Hair dryer|Bathtub|Shampoo|
      Essentials|Bidet|Conditioner|Body soap|Baby bath
      '),
32                  'bath-products'] = 1
33     listings.loc[listings['amenities'].str.contains('
      Bluetooth sound system|Ethernet connection|
      Heating|Pocket wifi|Cable TV|Wifi'),
34                  'electric-system'] = 1
35     listings.loc[listings['amenities'].str.contains('
      Breakfast'),
36                  'food-services'] = 1
37     listings.loc[listings['amenities'].str.contains('
      Outdoor furniture|Dining table|Hangers|High
      chair|Crib|Clothing storage: wardrobe|Dedicated

```

```

workspace|Drying rack for clothing|Bed linens|
      Extra pillows and blankets'),
38                  'house-furniture'] = 1
39     listings.loc[listings['amenities'].str.contains('
      Cleaning before checkout|Luggage dropoff allowed
      |Long term stays allowed'),
40                  'house-rules'] = 1
41     listings.loc[listings['amenities'].str.contains('
      Oven|Hot water kettle|Kitchen|Cooking basics|
      Microwave|Fire pit|Dishes and silverware|
      Barbecue utensils|Cleaning products|Baking sheet
      |Free washer|Free dryer|Iron|Dishwasher|Freezer|
      Coffee maker|Refrigerator|Toaster|dinnerware|BBQ
      grill|Stove|Wine glasses'),
42                  'kitchen-appliances'] =
43     1
44     listings.loc[listings['amenities'].str.contains('
      Free parking on premises|Free street parking'),
45                  'parking'] = 1
46     listings.loc[listings['amenities'].str.contains('
      Board games|Indoor fireplace|Bikes|Shared patio
      or balcony|Private fenced garden or backyard|
      crib|books and toys|Outdoor dining area|Private
      gym in building|Piano|HDTV with Netflix|premium
      cable|standard cable'),
47                  'recreation'] = 1
48     listings.loc[listings['amenities'].str.contains('
      Fire extinguisher|Carbon monoxide alarm|Window
      guards|Fireplace guards|First aid kit|Baby
      monitor|Private entrance|Lockbox|Smoke alarm|
      Room-darkening shades|Baby safety gates'),
49                  'safety'] = 1
50     host_verification = listings['host_verifications'
      ].str.split(',', expand=True)
51
52     listings.loc[listings['host_verifications'].str.
      contains('email'),
53                  'host_email'] = 1
54     listings.loc[listings['host_verifications'].str.
      contains('phone'),
55                  'host_phone'] = 1
56     listings.loc[listings['host_verifications'].str.
      contains('work_email'),
57                  'host_work_email'] = 1
58
59     new_feature_cols = listings.iloc[:,75:].columns
60     listings[new_feature_cols] = listings[
      new_feature_cols].fillna(0)
61     listings = listings.merge(reviews, how='inner',
      left_on='id', right_on='listing_id')
62
63     return listings
64
65 def plot_nas(df: pd.DataFrame):
66     if df.isnull().sum().sum() != 0:
67         na_df = (df.isnull().sum() / len(df)) * 100
68         na_df = na_df.drop(na_df[na_df == 0].index).
      sort_values(ascending=False)
69         missing_data = pd.DataFrame({'NaN %': na_df
      })
70         missing_data.plot(kind = "bar")
71         plt.show()
72     else:
73         print('No NAs found')
74
75 def drop_useless_cols(listings):
76     not_needed_columns = [
77         'id', 'listing_url', 'scrape_id', 'last_scraped',
78         'source', 'name',
79         'picture_url', 'host_id', 'host_url', 'host_name',
80         'host_location',
81         'host_about', 'host_thumbnail_url', '
      host_picture_url', 'host_neighbourhood',

```

```

80     'neighbourhood', 'neighbourhood_group_cleansed',
81     'calendar_updated', 'first_review', 'last_review',
82     'license',
83     'calculated_host_listings_count', '
84     calculated_host_listings_count_entire_homes',
85     'calculated_host_listings_count_private_rooms',
86     'calculated_host_listings_count_shared_rooms',
87     'description', 'neighborhood_overview', '
88     host_verifications', 'host_since',
89     'bathrooms', 'bathrooms_text', 'amenities', '
90     availability_30',
91     'availability_60', 'availability_90', '
92     availability_365', 'calendar_last_scraped',
93     'number_of_reviews_ltm', 'number_of_reviews_l30d',
94     'host_has_profile_pic', 'property_type',
95     'minimum_minimum_nights', '
96     maximum_maximum_nights', 'minimum_nights_avg_ntm',
97     '
98     minimum_maximum_nights', '
99     maximum_minimum_nights', 'maximum_nights_avg_ntm',
100 ]
101 listings.drop(not_needed_columns, axis = 1,
102             inplace = True)
103 listings = listings.dropna()
104
105 imputation_cols = ['bedrooms', 'beds']
106 for i in imputation_cols:
107     listings.loc[listings.loc[:,i].isnull(),i] =
108         listings.loc[:,i].median()
109
110 return listings
111
112 def plot_price_box(listings):
113     # Step 1: Clean the 'price' column
114     listings['price'] = listings['price'].astype(str).
115         str.replace('$', '').str.replace(',', '')
116     listings['price'] = pd.to_numeric(listings['price']
117                                     ])
118
119     # Step 2: Clean the 'host_response_rate' and '
120     host_acceptance_rate' columns
121     listings['host_response_rate'] = listings["
122     host_response_rate"].str.replace("%","")
123     listings['host_response_rate'] = pd.to_numeric(
124         listings['host_response_rate'])
125     listings['host_acceptance_rate'] = listings["
126     host_acceptance_rate"].str.replace("%","")
127     listings['host_acceptance_rate'] = pd.to_numeric(
128         listings['host_acceptance_rate'])
129
130     # Step 3: Plot the initial price boxplot
131     plt.figure(figsize=(10, 5))
132     listings[['price']].plot(kind='box', title='Price
133     BoxPlot')
134     plt.ylim(0, 1600)
135     plt.show()
136
137     # Step 4: Filter the 'price' column data
138     listings_filtered = listings[(listings.price > 50)
139                                & (listings.price <= 300)]
140
141     # Step 5: Plot the filtered price boxplot
142     plt.figure(figsize=(10, 5))
143     listings_filtered[['price']].plot(kind='box',
144                                     title='Filtered Price BoxPlot')
145     plt.show()
146
147     listings_filtered.head()
148
149     # Plot the distribution of host_response_rate
150     plt.figure(figsize=(10, 6))
151     listings['host_response_rate'].plot(kind='hist',
152                                     bins=20, edgecolor='black', alpha=0.7)
153     plt.title('Distribution of Host Response Rate')
154     plt.xlabel('Host Response Rate (%)')
155     plt.ylabel('Frequency')
156     plt.grid(axis='y', linestyle='--', alpha=0.7)
157     plt.show()
158
159     return listings
160
161 def plot_review_cols(listings):
162     # Convert review scores to numeric and fill NaNs
163     with the mean value of each column
164     review_score_columns = [
165         'review_scores_accuracy',
166         'review_scores_checkin',
167         'review_scores_cleanliness',
168         'review_scores_communication',
169         'review_scores_location',
170         'review_scores_rating',
171         'review_scores_value'
172     ]
173
174     # Ensure all columns are numeric and fill NaNs
175     with the mean value
176     for col in review_score_columns:
177         listings[col] = pd.to_numeric(listings[col],
178                                     errors='coerce')
179         listings[col].fillna(listings[col].mean(),
180                             inplace=True)
181
182     # Plot the distribution of each review score
183     column
184     fig, axes = plt.subplots(3, 3, figsize=(18, 18))
185
186     # Define positions to leave a gap in the grid
187     positions = [1, 2, 3, 4, 5, 6, 8] # Subplot
188     positions in a 3x3 grid
189
190     for i, (col, pos) in enumerate(zip(
191         review_score_columns, positions)):
192         ax = plt.subplot(3, 3, pos)
193         listings[col].plot(kind='hist', bins=20,
194                             edgecolor='black', alpha=0.7, ax=ax)
195         ax.set_title(f'Distribution of {col}')
196         ax.set_xlabel(f'{col}')
197         ax.set_ylabel('Frequency')
198         ax.grid(axis='y', linestyle='--', alpha=0.7)
199
200     plt.tight_layout()
201     plt.show()
202
203 def plot_neighborhood(listings):
204     # Calculate the number of listings in each
205     neighborhood
206     neighbourhood_DF = listings.groupby('
207     neighbourhood_cleansed').host_response_time.
208         count().reset_index()
209     neighbourhood_DF = neighbourhood_DF.rename(columns
210         ={'host_response_time': 'Number_Of_Listings'})
211
212     # Sort by the number of listings
213     neighbourhood_DF = neighbourhood_DF.sort_values(by
214         ='Number_Of_Listings', ascending=False)
215
216     # Plot the bar chart
217     plt.figure(figsize=(18, 8))
218     plt.bar(neighbourhood_DF['neighbourhood_cleansed']
219           , neighbourhood_DF['Number_Of_Listings'])
220     plt.title('Dublin Neighborhood Frequency')
221     plt.xlabel('Neighborhood')
222     plt.ylabel('Number of Listings')
223     plt.xticks(rotation=90) # Rotate x-axis labels
224     for better readability
225     plt.show()

```

```

189
190 def statics_review_score(listings):
191     # Define the list of review score columns
192     review_score_columns = [
193         'review_scores_accuracy',
194         'review_scores_checkin',
195         'review_scores_cleanliness',
196         'review_scores_communication',
197         'review_scores_location',
198         'review_scores_rating',
199         'review_scores_value'
200     ]
201
202     # Calculate and print statistics for each review
203     # score column
204     for col in review_score_columns:
205         print(f"{col} statistics:")
206         print(f"Mean: {listings[col].mean()}")
207         print(f"Median: {listings[col].median()}")
208         print(f"Mode: {listings[col].mode()[0]}")
209         print(f"Min: {listings[col].min()}")
210         print(f"Max: {listings[col].max()}")
211         print("\n")
212
213 def plot_Kernel_Density_Estimate(listings):
214     plt.figure(figsize=(18, 18))
215
216     # Define the list of review score columns and
217     # their subplot positions
218     review_score_columns = [
219         'review_scores_accuracy',
220         'review_scores_checkin',
221         'review_scores_cleanliness',
222         'review_scores_communication',
223         'review_scores_location',
224         'review_scores_rating',
225         'review_scores_value'
226     ]
227
228     positions = [1, 2, 3, 4, 5, 6, 8] # Subplot
229     # positions in a 3x3 grid
230
231     for col, pos in zip(review_score_columns,
232                         positions):
233         plt.subplot(3, 3, pos)
234         density = sns.kdeplot(listings[col], color="b"
235                               , label=f'{col} - density')
236         x = density.lines[0].get_xdata()
237         y = density.lines[0].get_ydata()
238         mean = listings[col].mean()
239         height = np.interp(mean, x, y)
240         density.vlines(mean, 0, height, color='b', ls=
241                        ':')
242         density.fill_between(x, 0, y, facecolor='b',
243                              alpha=0.2)
244         plt.legend()
245
246     plt.tight_layout()
247     plt.show()
248
249 def minmax(X):
250     X_std = (X - X.min()) / (X.max() - X.min())
251     X_scaled = X_std * (X.max() - X.min()) + X.min()
252     return X_scaled
253
254 def scaling_data(listings):
255     scaling_data = ['host_response_rate', '
256                     host_acceptance_rate', 'bedrooms', 'beds',
257                     'host_listings_count', '
258                     host_total_listings_count',
259                     'latitude', 'longitude', '
260                     accommodates', 'price',
261                     'minimum_nights', 'maximum_nights',
262                     'number_of_reviews', 'num_cmt', '#avg_senti',
263
264                     'review_scores_rating', '
265                     review_scores_accuracy',
266                     'review_scores_cleanliness', '
267                     review_scores_checkin',
268                     'review_scores_communication', '
269                     review_scores_location',
270                     'review_scores_value', '
271                     reviews_per_month', 'bath-products', 'electric-
272                     system',
273                     'food-services', 'house-furniture', '
274                     house-rules',
275                     'kitchen-appliances', 'parking', '
276                     recreation', 'safety',
277                     'host_email', 'host_work_email']
278
279     for i in scaling_data:
280         listings[i] = minmax(listings[i])
281
282     listings.dropna(axis = 1, inplace = True)
283     label_encoder = preprocessing.LabelEncoder()
284     listings.host_response_time = label_encoder.
285         fit_transform(listings.host_response_time)
286     listings.host_is_superhost = label_encoder.
287         fit_transform(listings.host_is_superhost)
288     listings.host_identity_verified = label_encoder.
289         fit_transform(listings.host_identity_verified)
290     listings.instant_bookable = label_encoder.
291         fit_transform(listings.instant_bookable)
292     listings.room_type = label_encoder.
293         fit_transform(listings.room_type)
294     listings.neighbourhood_cleansed = label_encoder.
295         fit_transform(listings.neighbourhood_cleansed)
296     listings.has_availability = label_encoder.
297         fit_transform(listings.has_availability)
298
299     test_corr = listings.corr()
300     test_corr.to_csv("test_corr.csv")
301
302     return listings
303
304 def plot_confusion_matrix(y_test, y_pred, title):
305     cm = confusion_matrix(y_test, y_pred)
306     plt.figure(figsize=(10, 7))
307     sns.heatmap(cm, annot=True, fmt='d', cmap='Blues
308                 ')
309     plt.xlabel('Predicted')
310     plt.ylabel('Actual')
311     plt.title(title)
312     plt.show()
313
314 def ROCAUC_visualizer(model, bin_count, X, y):
315     y = y['rating_bin_ep'].astype(int)
316     x_train, x_test, y_train, y_test =
317         train_test_split(X, y, test_size=0.2,
318                           random_state=1)
319     y_train = y_train.astype(int)
320     y_test = y_test.astype(int)
321
322     visualizer = ROCAUC(model, classes=list(range(
323         bin_count)), macro=False, micro=False)
324     visualizer.fit(x_train, y_train)
325     visualizer.score(x_test, y_test)
326     visualizer.show()
327
328 def bin_column(listings, col, n_bins):
329     X = listings[
330         ['host_response_time', '
331         host_response_rate', 'host_acceptance_rate',
332         'bedrooms', 'beds', '
333         neighbourhood_cleansed',
334         'host_is_superhost', '
335         host_listings_count', 'host_total_listings_count
336         ',
337         'host_identity_verified', '

```



```

304     'room_type',
305     'accommodates', 'price', '
306     'minimum_nights', 'maximum_nights',
307     'bath-products', 'electric-system',
308     'food-services', 'house-furniture',
309     'house-rules',
310     'kitchen-appliances', 'parking',
311     'recreation', 'safety',
312     'host_email', 'host_work_email'] +
313     list(reviews.columns[2:])
314 ]
315
316 y = listings[[col]]
317 y = (y/y.max())*100
318
319 y = y.assign(
320     rating_bin_ep = pd.qcut(
321         y[col],
322         q=n_bins,
323         duplicates='drop',
324         labels=list(range(n_bins))
325     )
326 )
327
328 y.groupby('rating_bin_ep').min()
329 y.groupby('rating_bin_ep').max()
330
331 return X , y
332
333 def select_important_features(X, y, threshold=0.05):
334
335     correlations = X.corrwith(y).abs().sort_values(
336         ascending=False)
337     selected_features = correlations[correlations >
338         threshold].index
339
340     plt.figure(figsize=(10, 6))
341     sns.barplot(x=selected_features, y=correlations[
342         selected_features], palette="viridis")
343     plt.axhline(y=threshold, color='r', linestyle='
344         --', label=f'Threshold ({threshold})')
345     plt.xlabel('Features')
346     plt.ylabel('Correlation with target')
347     plt.title('Selected Features and Their
348         Correlations')
349     plt.xticks(rotation=90)
350     plt.legend()
351     plt.show()
352
353     return selected_features, correlations
354
355 def check_bins(y):
356     y1 = y['rating_bin_ep']
357     cnt_plt = sns.countplot(y1)
358     cnt_plt.bar_label(cnt_plt.containers[0])
359     plt.show()
360
361 def evaluate_logistic_regression(X, y, c_range,
362     degree_range=[1], bin_count=2):
363
364     y1 = y['rating_bin_ep']
365     for i in degree_range:
366         trans = PolynomialFeatures(degree=i)
367         x_poly = trans.fit_transform(X)
368         x_train, x_test, y_train, y_test =
369         train_test_split(x_poly, y1, test_size=0.2,
370             random_state=1)
371         mean_error = []
372         std_error = []
373         for c in c_range:
374             log_reg = LogisticRegression(C=c,
375                 random_state=0, solver='newton-cg', multi_class='
376                 multinomial')
377             log_reg.fit(x_train, y_train)
378
379             y_pred = log_reg.predict(x_test)
380
381             cnf_mtx = confusion_matrix(y_test,
382                 y_pred)
383             f1_score = metrics.f1_score(y_test,
384                 y_pred, average='weighted')
385
386             scores = cross_val_score(log_reg, x_test
387                 , y_test, cv=5, scoring='accuracy')
388             mean_error.append(np.array(scores).mean
389                 ())
390             std_error.append(np.array(scores).std())
391
392             print(" Logistic Regression")
393             print(" For Degree = ", i)
394             print(" For C = ", c)
395             print(" Confusion Matrix - \n", cnf_mtx)
396             print(' Train accuracy score: ', log_reg
397                 .score(x_train, y_train))
398             print(' Test accuracy score: ', log_reg.
399                 score(x_test, y_test))
400             print(" F1 Score = ", f1_score)
401             print(" Mean Squared Error = ",
402                 mean_squared_error(y_test, y_pred))
403             print(" Classification Report\n",
404                 classification_report(y_test, y_pred))
405             print("\n")
406
407             title = f"Logistic Regression Confusion
408                 Matrix\nFor Degree = {i} and C = {c}"
409             plot_confusion_matrix(y_test, y_pred,
410                 title)
411
412             plt.errorbar(c_range, mean_error, yerr=
413                 std_error, linewidth=3)
414             plt.xlabel('C', fontsize=25)
415             plt.ylabel('Accuracy Score', fontsize=25)
416             title_cv = f"Logistic Regression - Cross
417                 Validation \nFor Degree = {i}"
418             plt.title(title_cv, fontsize=25)
419             plt.show()
420
421             ROCAUC_visualizer(log_reg, bin_count, x_poly
422                 , y)
423
424 def evaluate_knn(X, y, nn_range, bin_count=2):
425     y1 = y['rating_bin_ep']
426     x_train, x_test, y_train, y_test =
427     train_test_split(X, y1, test_size=0.2,
428         random_state=1)
429     merr = []
430     serr = []
431
432     for nn in nn_range:
433         knn_model = KNeighborsClassifier(n_neighbors
434             =nn, weights='uniform')
435         knn_model.fit(x_train, y_train)
436         y_pred_nn = knn_model.predict(x_test)
437
438         cnf_mtx = confusion_matrix(y_test, y_pred_nn
439             )
440         f1_score = metrics.f1_score(y_test,
441             y_pred_nn, average='weighted')
442
443         scores_knn = cross_val_score(knn_model,
444             x_test, y_test, cv=5, scoring='accuracy')
445         merr.append(np.array(scores_knn).mean())
446         serr.append(np.array(scores_knn).std())
447
448         print(" K Neighbors Classifier")
449         print(" For NN = ", nn)
450         print(" Confusion Matrix - \n", cnf_mtx)
451         print(' Train accuracy score: ', knn_model.
452             score(x_train, y_train))

```

```

417     print(' Test accuracy score: ', knn_model.
score(x_test, y_test))
418     print(" F1 Score = ", f1_score)
419     print(" Mean Squared Error = ",
mean_squared_error(y_test, y_pred_nn))
420     print(" Classification Report\n",
classification_report(y_test, y_pred_nn))
421     print("\n")
422
423     title = f"k-NN Confusion Matrix\nFor NN = {
nn}"
424     plot_confusion_matrix(y_test, y_pred_nn,
title)
425
426     plt.errorbar(nn_range, merr, yerr=serr,
linewidth=3)
427     plt.xlabel('NN', fontsize=25)
428     plt.ylabel('Accuracy Score', fontsize=25)
429     title_cv = f"k-NN - Cross Validation"
430     plt.title(title_cv, fontsize=25)
431     plt.show()
432
433     ROCAUC_visualizer(knn_model, bin_count, X, y)
434
435 def evaluate_decision_tree(X, y, depth_range,
bin_count=2):
436     y1 = y['rating_bin_ep']
437     x_train, x_test, y_train, y_test =
train_test_split(X, y1, test_size=0.2,
random_state=1)
438     merr_dt = []
439     serr_dt = []
440
441     for depth in depth_range:
442         dt_model = DecisionTreeClassifier(max_depth=
depth, random_state=1)
443         dt_model.fit(x_train, y_train)
444         y_pred_dt = dt_model.predict(x_test)
445
446         cnf_mtx = confusion_matrix(y_test, y_pred_dt
)
447         f1_score = metrics.f1_score(y_test,
y_pred_dt, average='weighted')
448
449         scores_dt = cross_val_score(dt_model, x_test
, y_test, cv=5, scoring='accuracy')
450         merr_dt.append(np.array(scores_dt).mean())
451         serr_dt.append(np.array(scores_dt).std())
452
453         print(" Decision Tree Classifier")
454         print(" For Depth = ", depth)
455         print(" Confusion Matrix - \n", cnf_mtx)
456         print(' Train accuracy score: ', dt_model.
score(x_train, y_train))
457         print(' Test accuracy score: ', dt_model.
score(x_test, y_test))
458         print(" F1 Score = ", f1_score)
459         print(" Mean Squared Error = ",
mean_squared_error(y_test, y_pred_dt))
460         print(" Classification Report\n",
classification_report(y_test, y_pred_dt))
461         print("\n")
462
463         title = f"Decision Tree Confusion Matrix\
nFor Depth = {depth}"
464         plot_confusion_matrix(y_test, y_pred_dt,
title)
465
466         plt.errorbar(depth_range, merr_dt, yerr=serr_dt,
linewidth=3)
467         plt.xlabel('Depth', fontsize=25)
468         plt.ylabel('Accuracy Score', fontsize=25)
469         title_cv = f"Decision Tree - Cross Validation"
470         plt.title(title_cv, fontsize=25)
471         plt.show()
472
473         ROCAUC_visualizer(dt_model, bin_count, X, y)
474
475     def plot_kde_by_bin(X, y, feature, bin_column, bins
=3, ylim=(0, 19), xlim=(-0.25, 1)):
476         test = pd.concat([X, y], axis=1)
477         colors = ['r', 'g', 'b']
478         labels = [f'{feature}, bin={i}' for i in range(
bins)]
479
480         for i in range(bins):
481             sns.kdeplot(test.loc[test[bin_column] == i,
feature], color=colors[i], label=labels[i]).set(
ylim=ylim, xlim=xlim)
482             plt.legend()
483
484         plt.xlabel(feature)
485         plt.ylabel('Density')
486         plt.title(f'KDE of {feature} by {bin_column}')
487         plt.show()
488
489     def plot_super_host(X, y):
490         test = pd.concat([X, y], axis=1)
491         test = test.groupby(['host_is_superhost', '
rating_bin_ep']).size().unstack()
492
493         plt.figure(figsize=(18, 6))
494
495         plt.subplot(1, 2, 1)
496         for j in range(3):
497             plt.bar(j, test.loc[1, j], label=f'Bin {j}',
alpha=0.7)
498         plt.title("Host is a Super Host")
499         plt.legend()
500         plt.xlabel('Review Bin', fontweight='bold',
fontsize=15)
501         plt.ylabel('# Records', fontweight='bold',
fontsize=15)
502
503         plt.subplot(1, 2, 2)
504         for j in range(3):
505             plt.bar(j, test.loc[0, j], label=f'Bin {j}',
alpha=0.7)
506         plt.title("Host is not a Super Host")
507         plt.legend()
508         plt.xlabel('Review Bin', fontweight='bold',
fontsize=15)
509         plt.ylabel('# Records', fontweight='bold',
fontsize=15)
510
511         plt.show()
512
513     # load csv
514     listings = pd.read_csv("listings.csv")
515     reviews = pd.read_csv("reviews_final.csv")
516
517     # listings pre-processing
518     listings = processing_and_merge(listings, reviews)
519     plot_nas(listings)
520     listings = drop_useless_cols(listings)
521     listings = plot_price_box(listings)
522     plot_review_cols(listings)
523     plot_neighborhood(listings)
524     statics_review_score(listings)
525     plot_Kernel_Density_Estimate(listings)
526     listings = scaling_data(listings)
527
528     # model parmater
529     c_range = [0.001, 0.1, 1, 10, 100, 1000]
530     degree_range = [1, 2]
531     nn_range = [1, 3, 5, 7, 9, 11, 13, 15, 17, 19]
532     depth_range = [3, 5, 7, 9, 11, 13, 15, 17, 19]
533

```

```

534 # review_scores_checkin evaluating
535 X,y = bin_column(listings, 'review_scores_checkin',
536 2)
537 check_bins(y)
538 yy = listings['review_scores_checkin']
539 select_important_features(X, yy, threshold=0.01)
540 evaluate_logistic_regression(X, y, c_range,
541 degree_range, 2)
542 evaluate_knn(X, y, nn_range, 2)
543 evaluate_decision_tree(X, y, depth_range, 2)
544
545 # review_scores_communication evaluating
546 X,y = bin_column(listings, '
547 review_scores_communication', 2)
548 check_bins(y)
549 yy = listings['review_scores_communication']
550 select_important_features(X, yy, threshold=0.01)
551 evaluate_logistic_regression(X, y, c_range,
552 degree_range, 2)
553 evaluate_knn(X, y, nn_range, 2)
554 evaluate_decision_tree(X, y, depth_range, 2)
555
556 # review_scores_cleanliness evaluating
557 X,y = bin_column(listings, '
558 review_scores_cleanliness', 3)
559 check_bins(y)
560 yy = listings['review_scores_cleanliness']
561 select_important_features(X, yy, threshold=0.01)
562 evaluate_logistic_regression(X, y, c_range,
563 degree_range, 3)
564 evaluate_knn(X, y, nn_range, 3)
565 evaluate_decision_tree(X, y, depth_range, 3)
566
567 # review_scores_location evaluating
568 X,y = bin_column(listings, 'review_scores_location',
569 3)
570 check_bins(y)
571 yy = listings['review_scores_location']
572 select_important_features(X, yy, threshold=0.01)
573 evaluate_logistic_regression(X, y, c_range,
574 degree_range, 3)
575 evaluate_knn(X, y, nn_range, 3)
576 evaluate_decision_tree(X, y, depth_range, 3)
577
578 # review_scores_rating evaluating
579 X,y = bin_column(listings, 'review_scores_rating',
580 3)
581 check_bins(y)
582 yy = listings['review_scores_rating']
583 select_important_features(X, yy, threshold=0.01)
584 evaluate_logistic_regression(X, y, c_range,
585 degree_range, 3)
586 evaluate_knn(X, y, nn_range, 3)
587 evaluate_decision_tree(X, y, depth_range, 3)
588
589 # review_scores_accuracy evaluating
590 X,y = bin_column(listings, 'review_scores_accuracy',
591 3)
592 check_bins(y)
593 yy = listings['review_scores_accuracy']
594 select_important_features(X, yy, threshold=0.01)
595 evaluate_logistic_regression(X, y, c_range,
596 degree_range, 3)
597 evaluate_knn(X, y, nn_range, 3)

```

```

595 evaluate_decision_tree(X, y, depth_range, 3)
596
597 # kde term home and num cmt
598 plot_kde_by_bin(X, y, 'term_home', 'rating_bin_ep',
599 bins=3, ylim=(0, 19), xlim=(-0.25, 1))
600 plot_kde_by_bin(X, y, 'num_cmt', 'rating_bin_ep',
601 bins=3, ylim=(0, 0.03), xlim=(-200, 1600))
602
603 # superhost
604 plot_super_host(X, y)

```

Listing 2. Python code of supplement.py

```

1 import random
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.model_selection import train_test_split
9 from sklearn.dummy import DummyClassifier
10 from sklearn.metrics import confusion_matrix
11 from sklearn.metrics import classification_report
12
13 def binary_step(x, thresh=10):
14     return np.where(x<thresh, 0, 1)
15
16 def decision_boundary(model):
17     b = model.intercept_[0]
18     w1, w2 = model.coef_[0]
19
20     c = -b / w2
21     m = -w1 / w2
22
23     xd = np.linspace(X.x_1.min(), X.x_1.max())
24     yd = m * xd + c
25
26     return xd, yd
27
28 size = 250
29 # NO LINEAR CORELATION
30 X = pd.DataFrame()
31 X['x_1'] = np.random.rand(size, )
32 X['x_2'] = np.random.randint(0, 60, size=size)
33
34 y = binary_step(10 + np.random.normal(0.0,1.0,size))
35
36 X_train, X_test, y_train, y_test = train_test_split(
37     X, y, test_size=0.2,
38
39     random_state=42)
40
41 model_no_linear_corr = LogisticRegression().fit(
42     X_train, y_train)
43 ypred = model_no_linear_corr.predict(X)
44
45 print("Logistic Regression - No Colinearity")
46 print(f"Model Score: {model_no_linear_corr.score(X,
47 y)}")
48 print(confusion_matrix(y, model_no_linear_corr.
49 predict(X)))
50 print(classification_report(y, model_no_linear_corr.
51 predict(X)))
52
53 pred_actual_logi = pd.DataFrame({'x_1': X['x_1'],
54 'x_2': X['x_2'],
55 'Ground Truth': y,
56 'Prediction': ypred})
57
58 # sns.scatterplot(x=X.x_2, y=X.x_1, style=y, hue=y,
59 palette='deep')
60
61 xd, yd = decision_boundary(model_no_linear_corr)

```

```

55 for idx, val in enumerate(yd):
56     if val < 0:
57         yd[idx] = 0
58     elif val > 50:
59         yd[idx] = 50
60
61 sns.scatterplot(x=X.x_1, y=X.x_2, style=yd, hue=
62     yd, palette='deep')
63 plt.plot(xd, yd)
64
65 # DATA IMBALANCE
66 X = pd.DataFrame()
67 size = 2000
68 X['x_1'] = np.arange(0,2,0.05)
69 y = 10 * X.x_1 + np.random.normal(0.0,1.0,X.size)
70
71 for idx, val in enumerate(y):
72     if val < 2:
73         y[idx] = int(0)
74     else:
75         y[idx] = int(1)
76 y = y.astype('int')
77
78 model_data_imbalance = LogisticRegression().fit(X, y)
79
80 print("Logistic Regression - Data Imbalance")
81 print(f"Model Score: {model_data_imbalance.score(X,
82     y)}")
83 print(confusion_matrix(y, model_data_imbalance.
84     predict(X)))
85 print(classification_report(y, model_data_imbalance.
86     predict(X)))
87
88 sns.scatterplot(x=X.x_1, y=y, style=y, hue=y,
89     palette='deep')
90 plt.legend(loc='lower right')
91 plt.ylabel('Target')
92 plt.xlabel('Feature')
93
94 ypred = model_data_imbalance.predict(X)
95 sns.scatterplot(x=X.x_1, y=ypred, style=ypred, hue=
96     ypred, palette='deep')
97 plt.legend(loc='lower right')
98 ypred = model_data_imbalance.predict(X)
99 plt.ylabel('Predicted')
100 plt.xlabel('Feature')
101
102 pred_data_imbalance_logi = pd.DataFrame({'x_1': X['
103     x_1'],
104     'Ground Truth': y,
105     'Prediction': model_data_imbalance.
106     predict(X)})
107
108 # DUMMY CLASSIFIER
109 model_dummy = DummyClassifier(strategy='
110     most_frequent').fit(X, y)
111
112 print("Dummy Classifier - Data Imbalance")
113 print(f"Model Score: {model_dummy.score(X, y)}")
114 print(confusion_matrix(y, model_dummy.predict(X)))
115 print(classification_report(y, model_dummy.predict(X)
116     )))
117
118 pred_data_imbalance_dummy = pd.DataFrame({'x_1': X['
119     x_1'],
120     'Ground Truth': y,
121     'Prediction': model_dummy.predict(X)})
122
123 # HOLD OUT METHOD
124 import numpy as np
125 from sklearn.model_selection import train_test_split
126
127 from sklearn.linear_model import LinearRegression
128 from sklearn.metrics import mean_squared_error
129 import matplotlib.pyplot as plt
130 import pandas as pd
131
132 X = np.arange(0,1,0.05).reshape(-1, 1)
133 y = 10 * X + np.random.normal(0.0,1.0,X.size).
134     reshape(-1, 1)
135
136 intercept = []
137 slope = []
138 mean_error = []
139
140 for i in range(5):
141     Xtrain, Xtest, ytrain, ytest = train_test_split(
142         X,y,test_size=0.2)
143
144     model = LinearRegression().fit(Xtrain, ytrain)
145     ypred = model.predict(Xtest)
146
147     intercept.append(model.intercept_[0])
148     slope.append(model.coef_[0][0])
149     mean_error.append(mean_squared_error(ytest,
150         ypred))
151
152 print('Intercept: {:.2f}\nSlope: {:.2f}\nSquared
153     Error: {:.2f}'.format(model.intercept_[0],
154         model.coef_[0][0],
155         mean_squared_error(ytest, ypred)))
156
157 print('\n\n')
158
159 y_vals = model.intercept_ + X * model.coef_
160 plt.plot(X, y_vals, label='{:.2f}'.format(
161     mean_squared_error(ytest, ypred)))
162
163 vals = pd.DataFrame({
164     'intercept': intercept,
165     'slope': slope,
166     'mean_error': mean_error})
167
168 plt.scatter(X, y, c='black')
169 plt.xlabel('Input X')
170 plt.ylabel('Target y')
171 plt.legend(title="MSE", fancybox=True)
172 plt.show();
173
174 import numpy as np
175 arr = np.arange(3.0, 5.5, 0.05)
176
177 # LAGGED OUTPUT
178 import pandas as pd
179 data = pd.DataFrame({
180     'Date': ['2024-07-15', '2024-08-15', '2024-09-15',
181         '2024-10-15', '2024-11-15', '2024-12-15'],
182     'Product Usage': [
183         1000, 2000, 4000, 8000, 16000, 32000],
184     'tag_1': [None, 1000, 2000, 4000, 8000, 16000],
185     'tag_2': [None, None, 1000, 2000, 4000, 8000],
186     })
187 print(data)

```