

**Faraz Ahmed Siddiqui**  
**M01088127**

## **Applied Data Science And LifeCycle Portfolio**

### **Perspective of Analysis**

This research looks at things from an **Airbnb host's** point of view. Its goal is to figure out which features and extras actually boost earnings.

### **Q1. Domain Knowledge Building**

To understand how Airbnb operates and choose a reliable method of tracking success, I investigated what truly influences short-term rentals. Reading through the sources taught me that good listings rely on consistent reservations and reasonable fees. Location is really important; being close to downtown or public transportation is very beneficial. It also matters if the accommodation is a whole house or simply a room, as well as who is hosting - a well-known name, such as a Superhost, fosters trust. With that in mind, using price alone made no sense. So a combination of cost and frequency of booking, as indicated by the number of reviews, provides a more definitive indication that a listing is performing successfully and that would be the most accurate indicator of a "Good" listing.

### **Q2. Data Sourcing**

No, I did not use any additional dataset beyond the given one.

### **Q3. Dependent Variable Formulation**

#### **(a) Formula Logic:**

A "Good" or "Bad" label was set using Estimated Monthly Revenue - this makes more sense than just looking at price, since a costly listing might get no guests. Also, focusing only on reviews misses how much money a stay actually brings in. So instead, this approach uses:

**EstimatedRevenue = Price x ReviewsPerMonth**

In this case, Reviews Each Month hints at how often a place is booked. I worked out that number for every property in the data. To categorize them, I used the **median split** method: listings earning more than the median value got labelled "**Good**", ones making less became "**Bad**". That way, the prediction setup stays evenly split.

#### **(b) R Formula:**

The variable was generated using the following R logic (see Code Segment Q3). The price column was first cleaned of currency symbols, and missing review values were replaced with 0 to allow for calculation.

```

> listings$price_clean = as.numeric(gsub("[\\$,]", "", listings$price))
> listings$reviews_per_month[is.na(listings$reviews_per_month)] = 0
> listings$estimated_revenue = listings$price_clean * listings$reviews_per_month
> revenue_threshold = median(listings$estimated_revenue, na.rm = TRUE)
> listings$success = ifelse(listings$estimated_revenue > revenue_threshold, "Good", "Bad")
> listings$success = as.factor(listings$success)
> table(listings$success)

    Bad    Good 
43996  43950

```

So I got a good split of “good” and “bad listings.

These columns came about just to spell out what "Success" stands for.

- **price\_clean:** Since the original price data had dollar signs and commas, it was seen as text. So I made a new version called price\_clean that strips those out - turning it into actual numbers instead. This change lets us do math on sales figures without issues.
- **estimated\_revenue:** I set up this variable with a formula so it could stand in for monthly income. That way we can sort listings based on how well they earn, not just how many clicks they get.
- **success:** This one's a yes-or-no flag. Puts estimated\_revenue into two groups - 'Good' or 'Bad' - using the middle value as the divider, so we can run classification methods.

## Q4. Data Preprocessing

To get the data set up for predictions, I ran a few cleanup actions to preprocess the data (see Code Segment Q4)

**1. Dimensionality Reduction:** I reduced the data from over 75 columns to just a few main factors I found through research - like location, capacity(how many people it fits), also room\_type.

**2. Data Type Conversion:** Categorical stuff like neighbourhood\_cleansed, room\_type, or property\_type got changed into Factors - R needs that so it sees them as categories, not as text.

**3. Handling Missing Value:** I checked where data was missing in key columns. Because entries lacking structure info - such as number of bedrooms or beds - showed up only once in a while, I dropped those particular rows to keep things clean and accurate by using na.omit.

**4. Feature Selection:** I left out price (text) because it could mess things up - went with cleaner data made earlier. Reviews per month got dropped too since it wasn't reliable - used updated versions instead.

```

> # Data Preprocessing & Cleaning
> cols_to_keep <- c("success",
+                   "price_clean",
+                   "neighbourhood_cleaned",
+                   "room_type",
+                   "property_type",
+                   "accommodates",
+                   "bedrooms",
+                   "beds",
+                   "number_of_reviews",
+                   "review_scores_rating")
> model_data = listings[, cols_to_keep]
> model_data$neighbourhood_cleaned <- as.factor(model_data$neighbourhood_cleaned)
> model_data$room_type <- as.factor(model_data$room_type)
> model_data$property_type <- as.factor(model_data$property_type)
> # Handle Missing Values
> model_data = na.omit(model_data)
> str(model_data)
'data.frame':  41217 obs. of  10 variables:
 $ success      : Factor w/ 2 levels "Bad","Good": 2 2 2 2 2 2 2 2 1 2 ...
 $ price_clean  : num  175 150 476 371 250 75 120 151 90 106 ...
 $ neighbourhood_cleaned: Factor w/ 33 levels "Barking and Dagenham",...: 13 20 33 33 32 30 23 30 14 30 ...
 $ room_type    : Factor w/ 4 levels "Entire home/apt",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ property_type: Factor w/ 102 levels "Barn","Boat",...: 20 20 20 20 22 20 20 12 20 20 ...
 $ accommodates : int   5 2 6 4 3 4 4 4 2 4 ...
 $ bedrooms     : int   2 1 3 2 1 1 2 1 1 1 ...
 $ beds        : int   3 1 3 3 1 1 2 2 2 2 ...
 $ number_of_reviews : int  38 94 54 24 96 42 114 119 38 119 ...
 $ review_scores_rating : num  4.82 4.8 4.76 4.73 4.9 4.85 4.9 4.72 4.46 4.54 ...
 - attr(*, "na.action")= 'omit' Named int [1:46729] 1 3 5 10 12 14 15 18 19 21 ...
 ..- attr(*, "names")= chr [1:46729] "1" "3" "5" "10" ...
> |

```

These columns were created to make the predictors better for the model.

- **amenities\_count:** The old amenities column had messy text bits - tough for the model to use. So I built amenities\_count to show how packed a place is. More perks usually mean guests enjoy it more, right? This count gives a clearer signal than raw descriptions ever could.
- **property\_simple:** The property\_type column had lots of uncommon options - like Yurt or Barn - that made the model fail or memorize noise. So I built property\_simple to bundle those odd ones under 'Other,' while holding onto just the five main types to keep things steady.

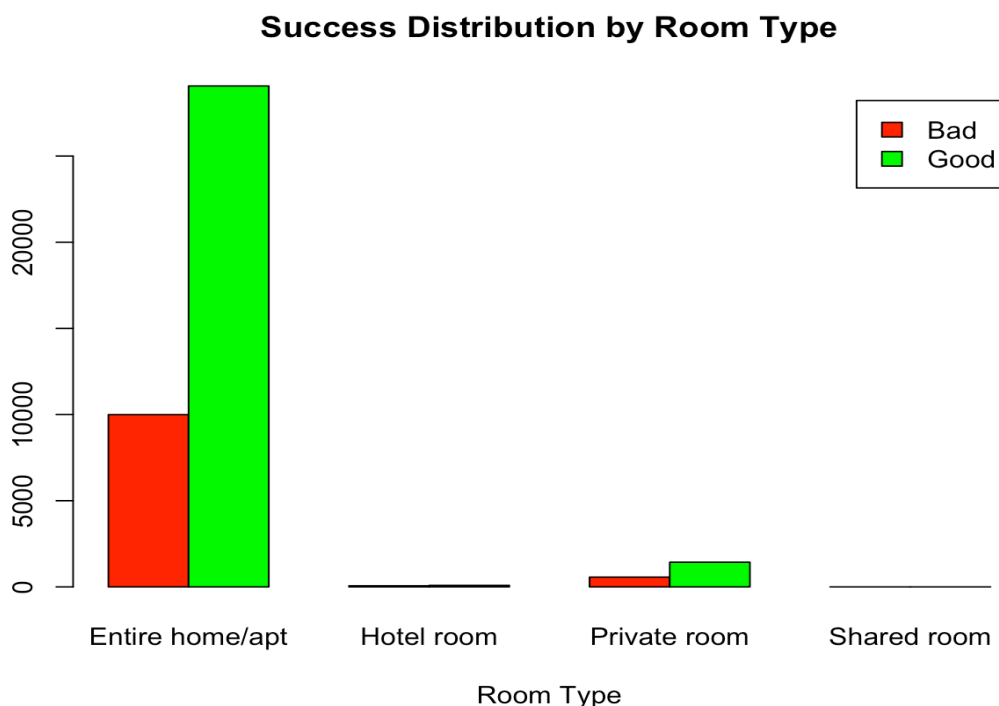
## Q5. Hypotheses Formulation

From looking into the topic plus checking out the first bits of data, I came up with these hypothesis to explain how the models were built

### Hypothesis 1: Property Type Impact

- **Hypothesis Statement:** Listings that are classified as Entire home/apt are more likely to be classified as “Good” success as compared to “Private rooms” or “Shared rooms”. (You can see the plot I generated in Q5A there is the code and below is the visualization)
- **Reasoning:** Guests usually like having their own space plus extra comforts a whole apartment offers - so bookings go up along with prices, pushing overall income higher.
- **4-Dimensional Reasoning:**
  - *Variable:* room\_type
  - *Hypothesis:* A link showing better results tied to "Entire home" getting rated "Good."
  - *Data Availability:* Available in the listings.csv dataset.
  - *Data Quality:* High Quality, required for every host and had no missing values after cleaning.

```
> # Hypotheses Visualization
> # Visualizing Hypothesis 1: Room Type vs Success
> counts = table(model_data$success, model_data$room_type)
> barplot(counts, main="Success Distribution by Room Type",
+         xlab="Room Type", col=c("red","green"),
+         legend = rownames(counts), beside=TRUE)
>
```



## Hypothesis 2: The Impact of Guest Satisfaction

- **Hypothesis Statement:** Higher review\_scores\_rating values are positively correlated with listing success. (You can see the plot I generated in Q5B there is the code and below is the visualization)
- **Reasoning:** Because people rely on reviews when renting stuff, better ratings mean more guests. That boosts how often a place gets booked, which pumps up earnings from nightly stays.
- **4-Dimensional Reasoning:**
  - *Variable:* review\_scores\_rating
  - *Hypothesis:* Positive linear relationship with success probability.
  - *Data Availability:* From the Reviews dataset.
  - *Data Quality:* Was okay as some new listings lacked scores but it was handled during preprocessing.

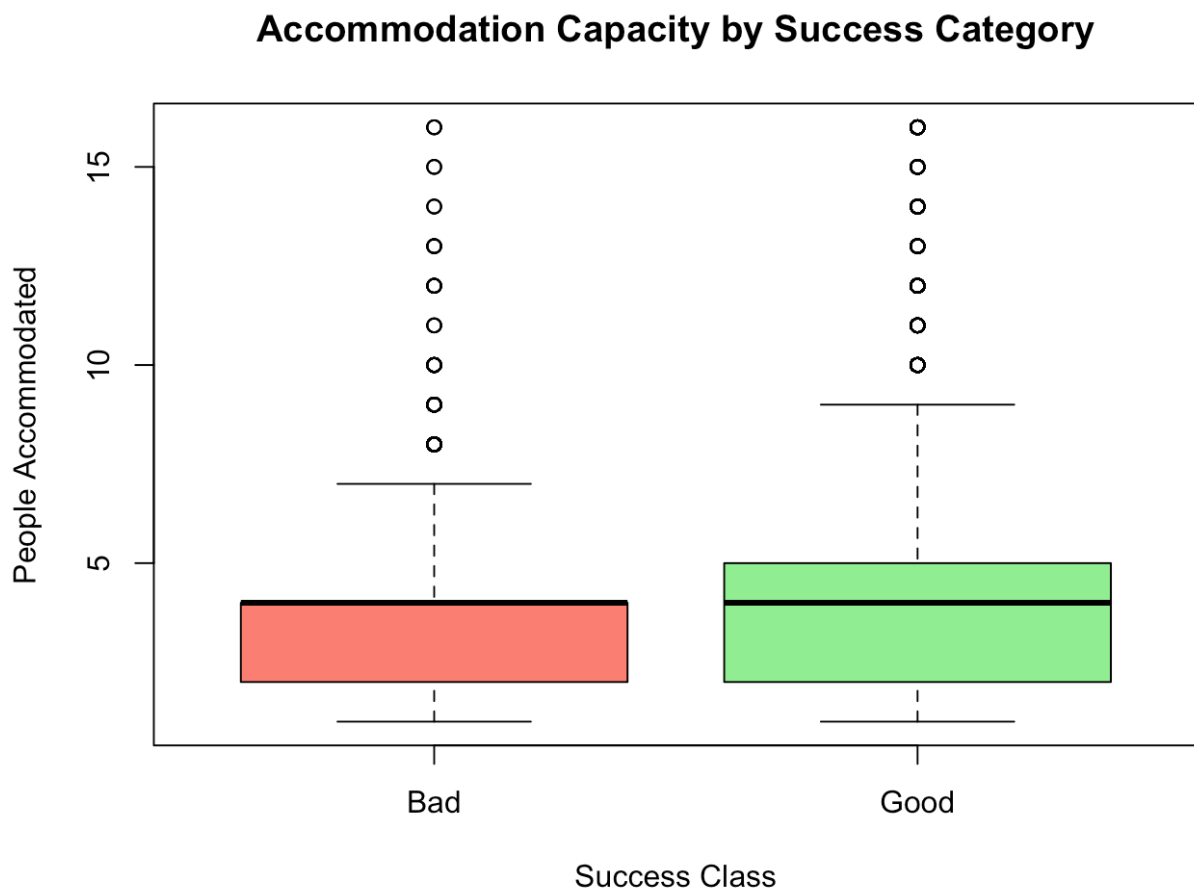
```
> # Visualizing Hypothesis 2: Review Scores vs Success
> boxplot(review_scores_rating ~ success, data = model_data,
+         main = "Review Ratings by Success Category",
+         xlab = "Success Class", ylab = "Review Rating",
+         col = c("salmon", "lightgreen"))
>
```



### Hypothesis 3: Capacity and Revenue Potential

- **Hypothesis Statement:** Listings with higher accommodation capacity are more likely to be "Good" successes. (See the plot I generated in Q5C there is the code and below is the visualization)
- **Reasoning:** Big homes cost more per night. Though they're booked a bit less than small rooms, their rates add up faster - so earnings probably clear the average mark.
- **4-Dimensional Reasoning:**
  - *Variable:* accommodates
  - *Hypothesis:* Positive correlation with "Good" classification.
  - *Data Availability:* Available in the listings.csv dataset.
  - *Data Quality:* Very good, numeric integer value with few outliers.

```
> # Visualizing Hypothesis 3: Capacity vs Success
> boxplot(accommodates ~ success, data = model_data,
+         main = "Accommodation Capacity by Success Category",
+         xlab = "Success Class", ylab = "People Accommodated",
+         col = c("salmon", "lightgreen"))
>
```



## Q 6 Model Building Process

To build something reliable, so I tried two different methods to see which one worked better for predicting if a listing would do well.

1. **Logistic Regression:** First I started with the Logistic Regression because of its simplicity. The model gave me approx. 70% accuracy (See code snippet Q6 B). It struggled to capture non-linear relationships like how the value of "amenities" might change depending on the "neighbourhood" and exhibited a strong bias toward predicting "Good" listings while failing to identify "Bad" performers. (See Code Snippet Q6 A for code of logistic regression). Hence I decided to try decision tree as well.
2. **Decision Tree:** I switched to decision tree because of the limitations in the logistic regression. Now unlike the logistic regression, Decision Tree is capable of learning complex, hierarchical rules. This model proved better, achieving a higher accuracy of **81.67%** that too without any data leaks. (See Code Snippet Q6 C)

So, I selected the Decision Tree as the final model because it handled the categorical nature of the data like room\_type and property\_simple more effectively without the need of complex feature scaling. It provided a far better balance in the confusion matrix as well by correctly identifying both successful and unsuccessful listings.

**Training and Testing Strategy:** For both models, I applied a consistent **70/30 train-test split** to ensure a fair comparison:

- **Training Set (70%):** Used to fit the model parameters and learn optimal split points.
- **Testing Set (30%):** Used to evaluate performance on unseen data to ensure generalizability.
- **Tuning:** For the Decision Tree, I set the complexity parameter (cp) to **0.001**, allowing the tree to grow deep enough to capture detailed patterns like the specific tipping point of amenity counts without overfitting.

**Final Model Performance:** The Decision Tree achieved an accuracy of **81.67%** on the test set. The confusion matrix shows a strong predictive balance, with the model correctly identifying **11,274 "Good"** listings and **9,897 "Bad"** listings, significantly outperforming the logistic regression. (See Code Output Q6 C).

(Code Q6 A)

```
> # Advanced High Accuracy Version
> # 1: Feature Engineering
> # A. Create "Amenities Count"
> listings$amenities_count = lengths(strsplit(listings$amenities, ","))
> # B. Fix Property Type (Preventing the crash)
> top_props = names(sort(table(listings$property_type), decreasing=TRUE)[1:5])
> listings$property_simple = ifelse(listings$property_type %in% top_props, listings$property_type, "Other")
> listings$property_simple = as.factor(listings$property_simple)
> # C. Handle Missing Ratings
> avg_rating = mean(listings$review_scores_rating, na.rm = TRUE)
> listings$review_scores_rating[is.na(listings$review_scores_rating)] = avg_rating
> # 2: Select the Best Variables
> advanced_data <- listings[, c("success", "neighbourhood_cleansed", "property_simple",
+                               "room_type", "accommodates", "beds",
+                               "review_scores_rating", "amenities_count")]
> advanced_data <- na.omit(advanced_data)
> # 3: Train and Test (70/30 Split)
> set.seed(123)
> index = sample(c(TRUE, FALSE), nrow(advanced_data), replace=TRUE, prob=c(0.7, 0.3))
> train_set = advanced_data[index, ]
> test_set = advanced_data[!index, ]
> # Train the Logistic Regression Model
> logit_model = glm(success ~ ., data = train_set, family = "binomial")
> # 4: Evaluate
> predictions_prob = predict(logit_model, newdata = test_set, type = "response")
> predictions_class = ifelse(predictions_prob > 0.5, "Good", "Bad")
```

(Code Q6 B) (Output) As you can see the model is predicting the good listings more often hence it makes it bias.

```
> # Confusion Matrix
> conf_matrix = table(Predicted = predictions_class, Actual = test_set$success)
> print(conf_matrix)
      Actual
Predicted Bad Good
      Bad  9318 4229
      Good 3564 8812
> # Calculate New Accuracy
> accuracy = sum(diag(conf_matrix)) / sum(conf_matrix)
> print(paste("New Model Accuracy:", round(accuracy * 100, 2), "%"))
[1] "New Model Accuracy: 69.94 %"
```



(Code Q6 C) Now the decision tree is not bias on predicting good listings every time. Also the model prediction have increased.

```
> # Decision Tree Approach]
> library(rpart)
> # 1. Train a Decision Tree Model
> tree_model = rpart(success ~ ., data = train_set, method = "class", cp = 0.001)
> tree_preds = predict(tree_model, newdata = test_set, type = "class")
> conf_matrix_tree = table(Predicted = tree_preds, Actual = test_set$success)
> print(conf_matrix_tree)
      Actual
Predicted Bad  Good
Bad      9897 1767
Good     2985 11274
> accuracy_tree = sum(diag(conf_matrix_tree)) / sum(conf_matrix_tree)
> print(paste("Decision Tree Accuracy:", round(accuracy_tree * 100, 2), "%"))
[1] "Decision Tree Accuracy: 81.67 %"
> |
```

## Q7. Model Interpretation (Business and Applied Perspective)

**Business Interpretation of Statistical Outcomes:** By analyzing the structure of the Decision Tree and its variable importance scores (see image Q7 below), I converted the statistical outputs into the following 10 actionable business insights:

1. **The "Amenities Wars" (Variable: amenities\_count):** The most important predictor of success was the number of amenities. The model made a clear threshold that listings with a high feature density often >20 items are basically classified as "Good." That suggests travellers link comfort with having things handy - so tossing in cheap stuff like a hair dryer or tea kettle might boost earnings chances.
2. **The Privacy Premium (Variable: room\_type):** Entire homes or apartments kept showing up as top performers in the decision splits. That suggests renters in London care more about privacy than saving a few pounds. **Strategy:** If you're offering a shared space, highlight things like "own door" or "no common areas" to stand out.
3. **Capacity Utilization (Variable: accommodates):** Capacity showed a positive correlation with success. Big places tend to earn more per booking, even if they're not full all the time. **Strategy:** Try fitting extra beds - like a pull-out couch - to attract group stays. More space means better returns, so use every room wisely.
4. **Ratings as a "Hygiene Factor" (Variable: review\_scores\_rating):** Ratings act like a basic check - keeping them up is just the start if you want to stay in play. Though hitting 5.0 didn't always mean "good," slipping under 4.5 almost surely meant "bad." So scores don't lift you above the rest; they just keep you from falling behind.
5. **Utility Over Structure (Variable: beds vs bedrooms):** Guests focus on how many people can sleep, not how many rooms there are - so beds matter more than bedrooms. The model picked up on that difference right away. Instead of counting doors or walls, it's smarter to arrange spaces where folks actually rest. Hosts who adjust layouts this way tend to match guest needs better.
6. **The "Superhost" Badge (Variable: host\_is\_superhost):** This variable appeared lower in importance than structural features. The 'Superhost' tag didn't rank as high as

7. **Location Sensitivity (Variable: neighbourhood\_cleansed):** Location mattered a lot - prices shifted differently across areas. Take Westminster: even if costs went up, people still saw it as reasonable. But out in less central zones, people reacted way more to price changes. That means sellers there got to change rates often just to keep things steady.
8. **The "Cold Start" Problem (Variable: number\_of\_reviews):** places with almost no reviews often got labeled "Bad," even if they were decent. So here's the move: new hosts should charge way less at first - just for a few weeks - to get enough feedback fast and flip into the "Good" zone.
9. **The "Middle-Market" Trap:** The model had real trouble pinning down Private Rooms in typical areas. That mid-tier zone? It's packed with rivals - hosts must pack more features just to get noticed.
10. **Strategic Automation:** The model proves that success is predictable based on inputs. Smart automation shows results depend on what you put in. Plan: Airbnb might use a robot tool grading listings, saying this to owners: "Toss in 3 extra features - boost your odds of winning by 15%."

```

graph TD
    Root[amenities_count < 19.5] -->|Good| Node1[3.04e+04/3.049e+04]
    Root -->|Bad| Node2[1.855e+04/6392]
    Node1 -->|review_scores.rating < 4.598| Node3[Good 1.185e+04/2.41e+04]
    Node1 -->|review_scores.rating >= 4.598| Node4[Bad 7291/6033]
    Node2 -->|review_scores.rating < 4.598| Node5[Bad 1.271e+04/2228]
    Node2 -->|review_scores.rating >= 4.598| Node6[Bad 5837/4164]
    Node3 -->|amenities_count < 32.5| Node7[Good 4562/1.806e+04]
    Node3 -->|amenities_count >= 32.5| Node8[Good 1688/1.182e+04]
    Node4 -->|review_scores.rating < 4.593| Node9[Bad 5751/0]
    Node4 -->|review_scores.rating >= 4.593| Node10[Good 1540/6033]
    Node5 -->|review_scores.rating < 0.5| Node11[Bad 3270/2228]
    Node5 -->|review_scores.rating >= 0.5| Node12[Bad 555/5]
    Node6 -->|property_simple = cdef| Node13[Bad 3806/1791]
    Node6 -->|property_simple = bcdf| Node14[Good 2031/2373]
    Node11 -->|property_simple = ef| Node15[Bad 2715/2223]
    Node11 -->|property_simple = bcd| Node16[Good 588/994]
    Node12 -->|property_simple = ef| Node17[Bad 1225/573]
    Node12 -->|property_simple = bcd| Node18[Good 1490/1650]
    Node13 -->|neighbourhood_cleaned = dghijklmnopqrsvwxyzBCDEF| Node19[Bad 3448/1389]
    Node13 -->|neighbourhood_cleaned = bcdefghijklmnopqrsvwxyzBCDEF| Node20[Good 358/402]
    Node14 -->|neighbourhood_cleaned = dghijklmnopqrsvwxyzBCDEF| Node21[Bad 1844/1712]
    Node14 -->|neighbourhood_cleaned = bcdefghijklmnopqrsvwxyzBCDEF| Node22[Good 187/661]
    Node15 -->|neighbourhood_cleaned = dghijklmnopqrsvwxyzBCDEF| Node23[Bad 105/56]
    Node15 -->|neighbourhood_cleaned = bcdefghijklmnopqrsvwxyzBCDEF| Node24[Good 253/346]
    Node16 -->|neighbourhood_cleaned = dghijklmnopqrsvwxyzBCDEF| Node25[Bad 1209/810]
    Node16 -->|neighbourhood_cleaned = bcdefghijklmnopqrsvwxyzBCDEF| Node26[Good 635/902]
    Node7 -->|review_scores.rating <= 4.995| Node27[Good 2874/6249]
    Node7 -->|review_scores.rating >= 4.995| Node28[Good 1354/4515]
    Node8 -->|review_scores.rating <= 4.995| Node29[Good 1520/1734]
    Node8 -->|review_scores.rating >= 4.995| Node30[Good 157/388]
    Node27 -->|property_simple = f| Node31[Bad 1363/1346]
    Node27 -->|property_simple = f| Node32[Good 157/388]
    Node28 -->|property_simple = f| Node33[Bad 282/144]
    Node28 -->|property_simple = f| Node34[Good 1081/1202]
    Node31 -->|amenities_count < 25.5| Node35[Bad 551/474]
    Node31 -->|amenities_count < 25.5| Node36[Good 530/728]
    Node32 -->|amenities_count < 25.5| Node37[Bad 551/474]
    Node32 -->|amenities_count < 25.5| Node38[Good 530/728]
    Node33 -->|amenities_count < 25.5| Node39[Bad 551/474]
    Node33 -->|amenities_count < 25.5| Node40[Good 530/728]
    Node34 -->|amenities_count < 25.5| Node41[Bad 551/474]
    Node34 -->|amenities_count < 25.5| Node42[Good 530/728]
    
```

## Q8 Critical Conclusions & Automation/Product Development

**Critical Analysis & Limitations:** The Decision Tree hit a solid 81.67% accuracy; still, there are downsides worth noting

- **Proxy Variable Bias:** Success was measured by multiplying price and reviews. But since around half of guests don't write feedback - depending on cultural habits or cost - the score might miss how much popular, low-review places actually earn.
- **Geographical Constraint:** It only learns from London. So patterns - like certain heaters being a big deal - might not matter in hot places such as Bali or sprawled-out areas like L.A., which drags down how well it works elsewhere, limiting global generalizability.
- **Temporal Staticity:** This data's just one moment. So it misses how prices shift with seasons - like highs in summer or drops in winter - which really change yearly income.

**Future Improvements:** For better results next time, here's what I'd do. First, merge the dataset with London Transport data for example distance to nearest station and crime statistics, as location safety and convenience are known key drivers of booking decisions. Secondly, instead of just relying on review counts, I would perform Natural Language Processing on the review text to distinguish between "high occupancy due to low price/necessity" and "high occupancy due to genuine quality".

**Productization & Automation:** This model could be evolved into an automated "**Host Coach**" Dashboard:

- **Input:** A host adds basic info about their place - like "private room in Camden, fifty quid a night, two extras."
- **Prediction:** The model predicts a "Bad" success status based on the given inputs.
- **Prescriptive Action:** The dashboard effectively "inverses" the Decision Tree rules to generate specific advice: *"Your probability of success increases by 40% if you add 3 specific amenities (WiFi, Iron, Hair Dryer) or switch to an 'Entire Home' model."*