

Predicting Airbnb Booking Popularity in Washington, DC.

Project Overview

This project explores a dataset of Airbnb listings in Washington, DC, to predict whether a listing has a "high booking rate," indicating frequent bookings.

Objectives

- **Exploration & Visualization:** Explore the data to understand structure, discover patterns, and handle missing values.
- **Model Development:** Build linear and logistic regression models to predict the target variable, high_booking_rate.
- **Interpretation & Insights:** Identify key features influencing booking popularity.
- **Model Evaluation:** Validate model performance on training and testing sets using relevant classification metrics.

Dataset Description

The dataset includes 4,936 Airbnb listings with property details, location, and booking information. The goal is to predict high_booking_rate, showing the likelihood of a listing being frequently booked.

Approach

- **Exploration & Visualization:** Identify key trends using Python.
- **Modeling:** Develop regression models, applying feature engineering to improve accuracy.
- **Validation:** Evaluate models based on metrics like accuracy, precision, recall, and AUC.

This project enhances my skills in data exploration, modeling, and predictive analysis within the short-term rental market context.

```
airbnb <- read_csv("dc_airbnb_listings.csv") # Load the dataset in R

## Rows: 4936 Columns: 12
## --- Column specification
## Delimiter: ";"
## chr (7): name, bed_type, cancellation_policy, cleaning_fee, price, property...
## dbl (5): accommodations, bedrooms, beds, host_total_listings_count, high_booki...
## Use 'spec()' to retrieve the full column specification for this data.
## Specify the column types or set 'show_col_types = FALSE' to quiet this message.

names(airbnb) # Variables used in dataset

## [1] "name" "accommodates"
## [2] "bed_type" "bedrooms"
## [3] "beds" "cancellation_policy"
## [7] "cleaning_fee" "host_total_listings_count"
## [8] "price" "property_type"
## [9] "room_type" "high_booking_rate"

age_mean <- airbnb %>%
  summarise(mean_accommodates = mean(accommodates))
```

1: EDA and Data Cleaning

a) To prepare the dataset for analysis, the following cleaning and preprocessing steps were performed:

- **Cancellation Policy Grouping:** Combined the "strict" and "super_strict_30" categories into a single "strict" category to simplify analysis.
- **Convert to Numerics:** Transformed cleaning_fee and price into numeric values.
- **Handling Missing Values (NAs):**
 - Replaced missing values in cleaning_fee and price with 0.
 - Imputed missing values in other numeric variables with their mean.

```
data.airbnb_clean <- airbnb %>%
  mutate(cancellation_policy = ifelse(cancellation_policy %in% c("strict", "super_strict_30"), "strict", cancella
tion_policy),
  cleaning_fee = as.numeric(gsub("[^0-9]", "", cleaning_fee)),
  price = as.numeric(gsub("[^0-9]", "", price)),
  cleaning_fee = ifelse(is.na(cleaning_fee), 0, cleaning_fee),
  price = ifelse(is.na(price), 0, price)
)

num_cols <- c("accommodates", "bedrooms", "beds", "host_total_listings_count")

data.airbnb_clean <- data.airbnb_clean %>%
  mutate_at(num_cols, ~replace_na(., mean(., na.rm = TRUE)))
```

b) Feature Engineering: Creating New Variables

To enhance the dataset, the following new variables were created:

- **price_per_person:** Calculated as the nightly price divided by the number of accommodates, representing the cost per person per night.
- **has_cleaning_fee:** A binary variable indicating whether a cleaning fee exists. Coded as "yes" if cleaning_fee > 0, "no" otherwise.
- **bed_category:** A categorical variable with values "bed" if bed_type is "Real bed" and "other" for all other bed types.
- **property_category:** Categorized based on the property_type:
 - "apartment" for "Apartment", " serviced apartment", or "loft".
 - "hotel" for "Bed & breakfast", "boutique hotel", or "hotel".
 - "condo" for "Condominium" or "condominium".
 - "house" for "Bungalow" or "house".
 - "other" for any other property types. This variable was converted to a factor for analysis.

- **ppp_ind:** A binary indicator where 1 denotes that price_per_person is greater than the median price_per_person within its property_category, and 0 otherwise.

```
data1 <- data.airbnb_clean %>%
  mutate(
    price_per_person = price/accommodates,
    has_cleaning_fee = ifelse(cleaning_fee > 0, "yes", "no"),
    bed_category = ifelse(bed_type == "Real Bed", "bed", "other"),
    property_type_new = case_when(
      property_type %in% c("Apartment", " serviced apartment", "loft") ~ "apartment",
      property_type %in% c("Bed & breakfast", "boutique hotel", "hotel") ~ "hotel",
      property_type %in% c("Condominium", "condominium") ~ "condo",
      property_type %in% c("Bungalow", "house") ~ "house",
      TRUE ~ "other"
    ) %>%
    group_by(property_category) %>%
    mutate(
      ppp_ind = ifelse(price_per_person > median(price_per_person, 1, 0)
    )
  ) %>%
  ungroup() %>%
  mutate(
    property_category = factor(property_category, levels = c("apartment", "hotel", "condo", "house", "other"))
  )
```

c) Converting the remaining character variables to factors:

```
## bed_type
## cancellation_policy
## room_type

data2 <- data1 %>%
  mutate(
    bed_type = factor(bed_type),
    cancellation_policy = factor(cancellation_policy),
    room_type = factor(room_type)
  )
```

d) Visualizing Price Per Person by Booking Rate

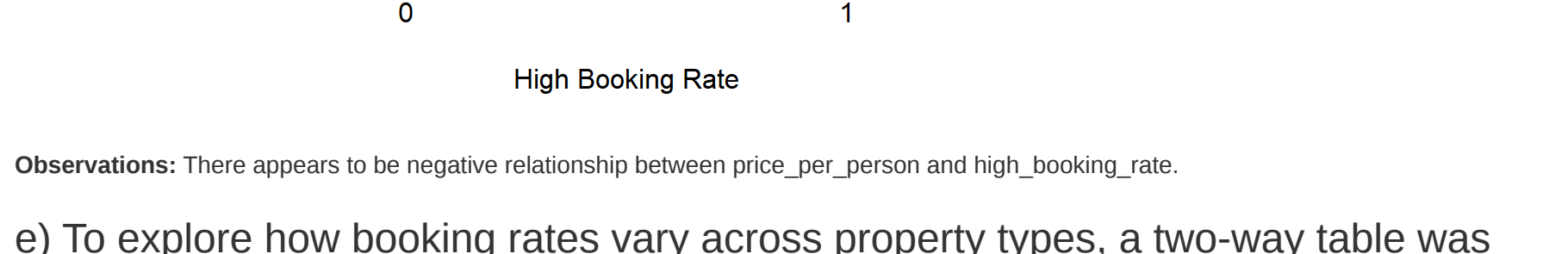
To examine the relationship between the nightly price per person and the booking rate, boxplots were constructed for both price_per_person and log(price_per_person), grouped by high_booking_rate.

- **Boxplots of price_per_person:** The distribution of price per person is visualized across different values of high_booking_rate.
- **Boxplots of log(price_per_person):** A logarithmic transformation is applied to price_per_person to better visualize the distribution and reduce the effect of outliers.

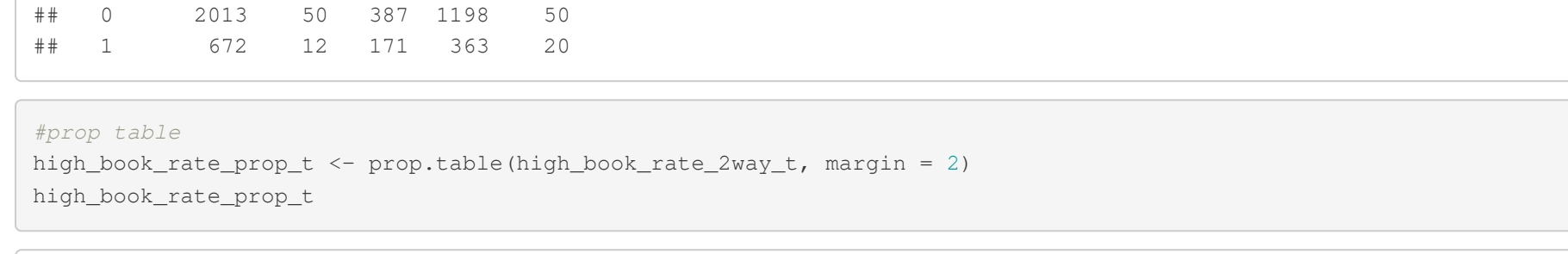
Observations:

```
library(ggplot2)

# Boxplot of price_per_person
ggplot(data1, aes(x = high_booking_rate, y = price_per_person)) +
  geom_boxplot() +
  labs(title = "Boxplot of price_per_person by high_booking_rate",
       xlab = "High Booking Rate",
       ylab = "Price Per Person",
       col = "bed")
```



```
# Boxplot of log(price_per_person)
ggplot(data1, aes(x = high_booking_rate, y = log(price_per_person))) +
  geom_boxplot() +
  labs(title = "Boxplot of log(price_per_person) by high_booking_rate",
       xlab = "High Booking Rate",
       ylab = "Log(Price Per Person)",
       col = "bed")
```



Observations: There appears to be a negative relationship between price_per_person and high_booking_rate.

e) To explore how booking rates vary across property types, a two-way table was constructed between high_booking_rate and property_category.

```
## two way table
high_book_rate_2way_t <- table(data2$high_booking_rate, data2$property_category)
high_book_rate_2way_t

##
##      apartment hotel condo house other
## 0      2013      50      387      198      50
## 1       672      12       171      363      20
```

```
## prop table
high_book_rate_prop_t <- prop.table(high_book_rate_2way_t, margin = 2)
high_book_rate_prop_t

##
##      apartment hotel condo house other
## 0 0.2492297 0.0045316 0.0395484 0.1764500 0.1242897
## 1 0.2502793 0.1935484 0.3064516 0.2326320 0.2857143
```

2: Linear Regression

a) Training a linear regression to predict high_booking_rate using the variables listed below.

- cancellation_policy
- cleaning_fee
- price_per_person
- ppp_ind
- has_cleaning_fee
- accommodations
- bed_category
- bedrooms
- beds
- host_total_listings_count
- property_category

```
model <- lm(high_booking_rate ~ cancellation_policy + cleaning_fee + price_per_person + ppp_ind + has_cleaning_f
ee + accommodations + bed_category + bedrooms + beds + host_total_listings_count + property_category, data = data2)
summary(model) # squared

## [1] 0.7898232
```

Model 1 has an R² of 0.7898.

b) Given a set of listing characteristics, the goal is to predict the high_booking_rate. The prediction is based on the model built earlier in the analysis.

New Listing:

- cancellation_policy = super_strict_30
- cleaning_fee = \$30
- price = \$200
- accommodations = 4
- bed_type = Real Bed
- bedrooms = 3
- beds = 4
- host_total_listings_count = 1
- property_type = townhouse

```
new_listing <- data.frame(cancellation_policy = "strict",
  cleaning_fee = 30, price_per_person = 200/4,
  ppp_ind = 1, has_cleaning_fee = "yes",
  accommodations = 4,
  bed_category = "other",
  bedrooms = 3,
  beds = 4,
  host_total_listings_count = 1,
  property_category = "condo"
)

high_booking_rate_pred <- predict(model, newdata = new_listing)
high_booking_rate_pred

## [1] 0.2602891
```

The high_booking_rate prediction value for the new listing, was close to 0 than 8 to 1, which suggests that our model has predicted that this new listing will not have a high booking rate.

3: Logistic Regression

a) Training a Logistic Regression model using the same variables as in Linear Regression model:

```
model_logit <- glm(high_booking_rate ~ cancellation_policy + cleaning_fee + price_per_person + ppp_ind + has_cle
aning_fee + accommodations + bed_category + bedrooms + beds + host_total_listings_count + property_category, data =
data2, family = binomial)

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

summary(model_logit)
```

```
## Call:
## glm(formula = high_booking_rate ~ cancellation_policy + cleaning_fee +
## price_per_person + ppp_ind + has_cleaning_fee + accommodations +
## bed_category + bedrooms + beds + host_total_listings_count +
## property_category, family = binomial, data = data2)
## Coefficients:
## (Intercept) Estimate Std. Error z value Pr(>|z|)
## cancellation_policystrict 0.510396 0.093378 5.462 6.42e-08 ***
## cancellation_policyother -0.444880 0.098816 4.444 3.42e-06 ***
## cleaning_fee -0.008897 0.002389 -3.708 0.00017 ***
## price_per_person -0.019105 0.002891 -6.608 3.91e-11 ***
## ppp_ind -0.023136 0.108811 -0.489 0.6247
## has_cleaning_feeYES 0.841786 0.074002 8.848 6.39e-15 ***
## accommodations 0.040726 0.035167 1.158 0.2468
## bed_categoryother -0.246419 0.026151 -9.375 0.00012 ***
## bedrooms -0.399557 0.069183 -5.775 7.69e-09 ***
## beds 0.223225 0.020891 4.388 3.16e-05 ***
## host_total_listings_count -0.007824 0.003384 -2.449 0.0143
## property_categoryhotel 0.100804 0.352788 0.286 0.7751
## property_categorycondo 0.349404 0.193860 1.814 0.0741
## property_categoryhouse -0.052321 0.081447 -0.627 0.5307
## property_categoryother -0.140461 0.291836 -0.481 0.6303
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## (Dispersion parameter for binomial family taken to be 1)
## Null deviance: 5560.1 on 4935 degrees of freedom
## Residual deviance: 5082.6 on 4920 degrees of freedom
## AIC: 5114.6
## Number of Fisher Scoring iterations: 8
```

The AIC of this model is 5114.6.

b) The coefficient for price per person is -0.019105, which means that for a one-unit increase in price per person, the log-odds of the high_booking_rate is would decrease by 0.019105, if all other variables in the model remain constant.

c) The coefficient for the condo property category is 0.349404, which means that the log odds of a condo listing having a high booking rate is expected to be 0.349404 higher than the log odds of the base property_category, if all else remain constant.

d) To estimate the likelihood that the given listing has high_booking_rate = 1, we use the logistic regression model and compute the predicted probability.

```
prob_pred <- predict(model_logit, newdata=new_listing, type = "response")
print(prob_pred)

##
## [1] 0.2743069
```

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069

[1] 0.2743069