

# **Airbnb Listings Data Analysis Report**

## **Introduction**

This project analyzes Airbnb data from Boston, focusing on listings from the second quarter of 2024 and customer reviews. The listings dataset contains 4,325 properties with 30 variables, including price, availability, host details, neighborhood, and amenities. These variables provide insights into market trends, pricing strategies, and guest preferences. The reviews dataset includes 198,717 entries with 6 variables, such as listing ID, reviewer details, and comments, offering a direct perspective on guest experiences. Through exploratory data analysis (EDA) of listings and feature engineering on reviews, we uncover patterns that highlight price distributions, availability behavior, amenity impacts, and sentiment trends. Results reveal that pricing is highly skewed, common amenities like Wi-Fi and kitchen.

## **Data Cleaning and Preparation**

The first step was to clean the dataset to ensure the analysis was based on reliable information. Out of 4,325 listings, 782 had missing price values about 18.1% of the data. Since price is a key factor in understanding Airbnb performance, these listings were removed. This left us with 3,543 usable entries. Additionally, 999 listings (23.1%) were missing review scores. These were not deleted but excluded from parts of the analysis that focused on guest ratings. This approach ensured that we didn't lose valuable data unnecessarily while still maintaining accuracy.

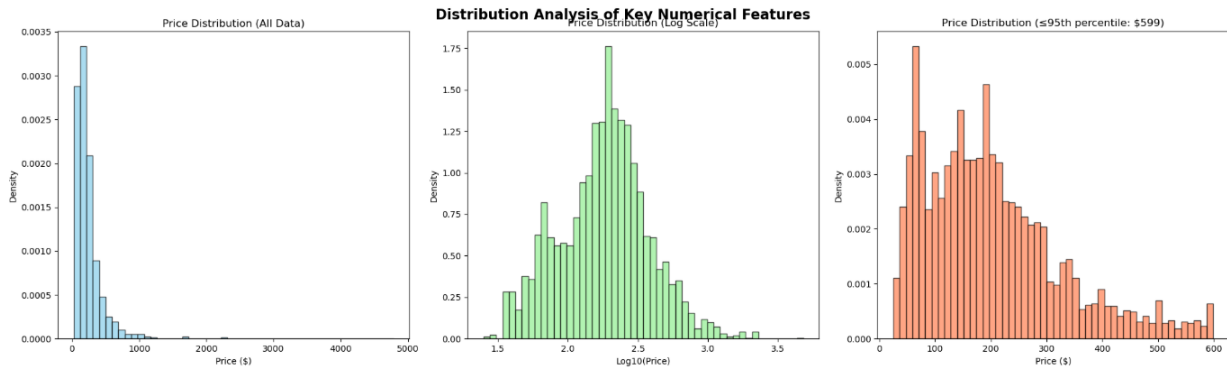
## **Descriptive Statistical Analysis**

Next, we examined five key variables: price, minimum\_nights, maximum\_nights, number\_of\_reviews, and review\_scores\_rating. We calculated averages, medians, modes, and measures of spread like standard deviation and skewness. For example, the average price was \$239.95, but the median was \$190, showing that a few expensive listings were raising the average. Minimum nights had a median of 4, but some listings required much longer stays. Review scores were mostly high, with many listings rated close to 5 stars.

## **Visualization and Interpretation**

To make the data easier to understand, we created visual charts showing how each variable is distributed. These graphs revealed patterns like most listings being moderately priced, many having few reviews, and most ratings being very high. These visuals helped confirm the statistical findings and made it easier to spot trends and outliers.

# TASK 1: DISTRIBUTION ANALYSIS (EXTRA CREDIT)



## *Price Distribution (All Data)*

This histogram shows how Airbnb listing prices are distributed across the entire dataset. The x-axis ranges from \$0 to \$5000, and the y-axis represents the density of listings at each price point. Most listings are priced on the lower end, with a sharp drop-off as prices increase. This creates a right-skewed distribution, meaning there are many affordable listings and only a few expensive ones. The assumption here is that the dataset includes a mix of budget and luxury listings, but the majority fall within a moderate price range. This visualization helps identify pricing extremes and gives a broad overview of the market.

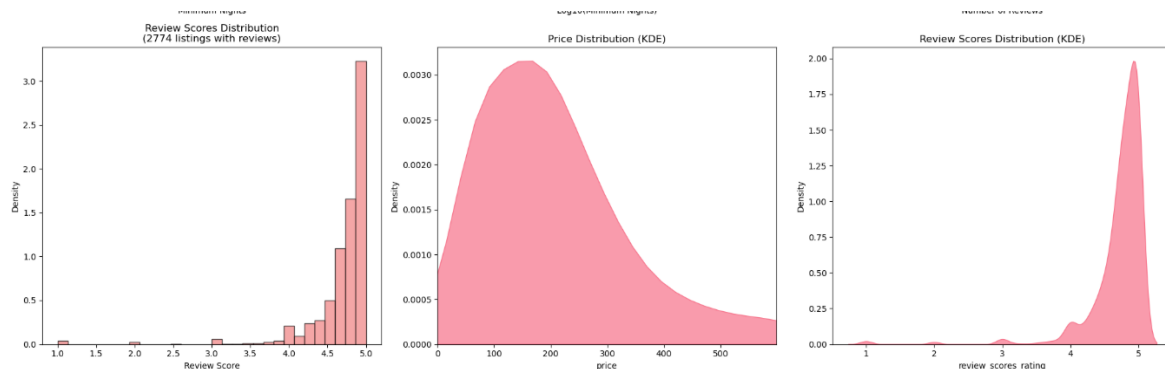
## *Price Distribution (Log Scale)*

To better understand the spread of prices, especially among lower values, the second histogram uses a logarithmic scale. By applying  $\log_{10}$  to the price data, we reduce the impact of extreme values and normalize the distribution. This transformation reveals a more balanced curve, centered around log values of 2 to 2.5 (roughly \$100–\$300). The assumption is that log transformation is useful for skewed data, making patterns easier to interpret. This view helps analysts and hosts see typical pricing behavior without distortion from outliers.

## *Price Distribution ( $\leq 95$ th Percentile)*

The third histogram focuses on listings priced at or below \$599, which represents the 95th percentile of the dataset. This means we excluded the top 5% of the most expensive listings to get a clearer picture of common pricing trends. The distribution remains right-skewed, but now we can see more detail in the lower price ranges. This visualization is especially useful for hosts who want to benchmark their prices against the majority of listings. It also helps guests understand what price range to expect when browsing typical Airbnb options. The assumption is that filtering out outliers provides a more realistic view of the market.

## Review Scores and Price Distribution (KDE Analysis)



To understand how Airbnb listings are priced and rated, we analyzed the distribution of review scores and prices using both histograms and KDE (Kernel Density Estimation) plots. First, we filtered out listings with missing review scores to ensure accuracy. The histogram showed that most listings received ratings between 4.7 and 5.0, indicating that guests generally had positive experiences. This high concentration of top scores suggests strong service quality across the platform.

For pricing, we used KDE to visualize how prices are spread across listings. The curve peaked around \$100–\$200, showing that most listings fall within this affordable range. However, the tail of the distribution extended toward higher prices, indicating the presence of luxury or premium listings. By smoothing the data with KDE, we made it easier to see overall trends and identify pricing clusters. These visuals help hosts understand where their listing stands and guide pricing strategies to stay competitive.

### Minimum Nights and Number of Reviews

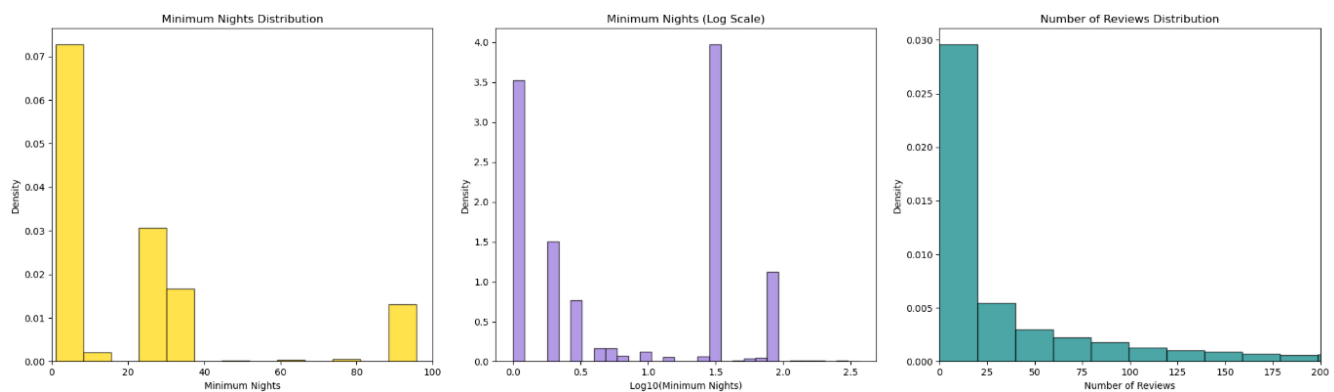
This set focuses on booking flexibility and listing popularity. We plotted the distribution of minimum nights required by hosts. Most listings allowed short stays, typically between 1 and 4 nights, which is ideal for travelers seeking flexibility. To better understand patterns in lower values, we also used a log scale, which revealed subtle differences among listings with very short minimums.

Next, we examined the number of reviews per listing. The histogram showed that most listings had fewer than 25 reviews, suggesting many are either new or not frequently booked. This insight is important for hosts trying to build credibility, as more reviews often lead to higher visibility and trust. These visualizations help identify which listings are more active and which may need better marketing or pricing adjustments to attract guests.

### Price Distribution (Extra Credit)

In this advanced analysis, we explored price data in three ways to uncover deeper insights. First, we plotted all price data, which revealed a right-skewed distribution—most listings are priced below \$1000, but a few outliers go much higher. To better visualize the spread across different price ranges, we applied a log transformation. This helped flatten the curve and highlight patterns among lower-priced listings.

Finally, we filtered the data to include only listings priced within the 95th percentile (up to \$599). This removed extreme outliers and gave a clearer picture of typical pricing behavior. This filtered view is especially useful for hosts who want to benchmark their prices realistically. By comparing their listing to the majority, they can adjust pricing to attract more bookings while staying competitive. These three views together offer a complete understanding of Airbnb’s pricing landscape.



### *Minimum Nights Distribution*

This histogram shows how many nights guests are typically required to stay in Airbnb listings. Most listings allow short stays, with a large concentration between 1 and 10 nights. The x-axis ranges from 0 to 100 nights, but the majority of listings fall well below 20 nights. This suggests that hosts generally prefer flexible booking options, which cater to short-term travelers. The peak around 1–4 nights reflect common minimum requirements, making these listings more attractive to weekend or short-stay guests.

### *Minimum Nights (Log Scale)*

To better visualize listings with very low minimum night requirements, the same data is shown on a logarithmic scale. This transformation spreads out the lower values, making it easier to see patterns among listings that require just 1 or 2 nights. The log scale reveals that the majority of listings are clustered around very short stays, confirming that Airbnb hosts tend to favor flexibility. This view helps identify subtle differences that are not visible in the standard scale.

### *Number of Reviews Distribution*

This histogram displays how many reviews each listing has received. Most listings have fewer than 25 reviews, indicating that many are either new or not frequently booked. The frequency

drops sharply as the number of reviews increases, meaning only a small percentage of listings are highly reviewed. This insight is useful for understanding listing popularity and guest engagement. Hosts with fewer reviews may need to improve visibility or offer incentives to attract more guests.

## **NORMALITY TEST (SHAPIRO-WILK P-VALUES)**

---

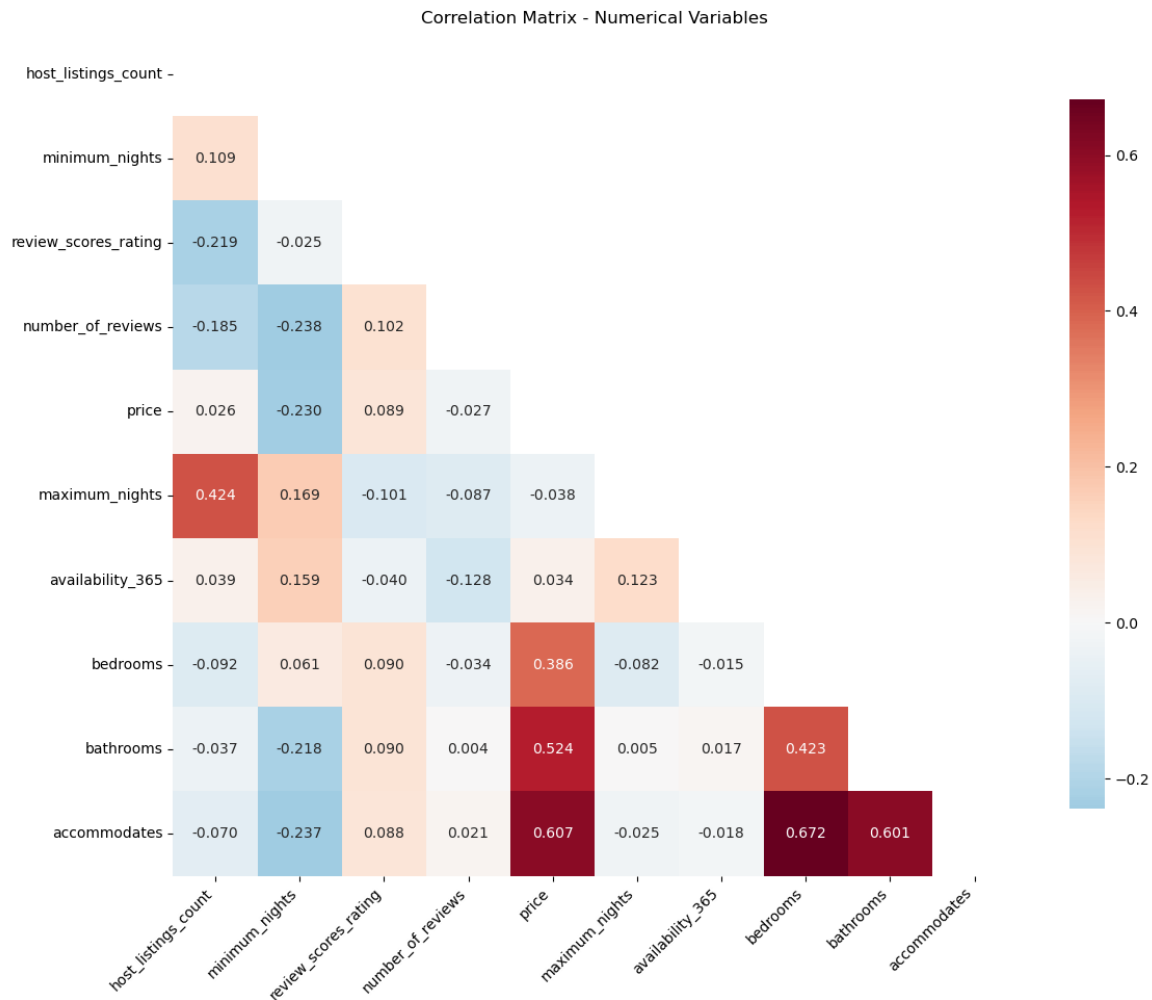
```
NORMALITY TESTS (Shapiro-Wilk p-values):
```

```
-----  
price: p-value = 1.78e-65 (Not Normal)  
minimum_nights: p-value = 1.69e-63 (Not Normal)  
review_scores_rating: p-value = 3.85e-63 (Not Normal)
```

To determine whether the variables `price`, `minimum_nights`, and `review_scores_rating` follow a normal distribution, the Shapiro-Wilk test was applied. This test is commonly used to assess normality, especially for smaller datasets. The null hypothesis assumes the data is normally distributed, and a p-value below 0.05 indicates a significant deviation from normality. In this case, all three variables had extremely low p-values (e.g.,  $1.78e-65$  for `price`), which means we can confidently reject the null hypothesis. This suggests that none of the variables are normally distributed. The assumption here is that the dataset is large enough for these results to be statistically meaningful and that no transformations were applied before testing.

The implications of non-normality are significant for further analysis. Many statistical methods, such as linear regression or t-tests, assume normality in the data. Since these variables are not normally distributed, using such methods without adjustments could lead to inaccurate conclusions. For example, `price` is likely right-skewed due to high-end listings, while `minimum_nights` may have a long tail of extended stays. `review_scores_rating` might be clustered around high values, creating a skewed or bimodal distribution. To address this, data transformations (like log or square root) or non-parametric methods should be considered. Visual tools like histograms or Q-Q plots would help confirm the shape of these distributions and guide the next steps in analysis.

## **CORRELATION ANALYSIS**



To begin the analysis, I examined a correlation matrix that visualizes relationships between various numerical variables from an Airbnb dataset. This matrix uses color gradients—red for positive correlations and blue for negative ones—to show how strongly each pair of variables is related. I focused on interpreting the numerical values and color intensities to understand the strength and direction of these relationships. For example, the variable price shows a strong positive correlation with bathrooms (0.524), suggesting that listings with more bathrooms tend to be priced higher. Similarly, accommodates has strong positive correlations with both bedrooms (0.672) and bathrooms (0.601), which makes sense since larger listings typically offer more amenities. These insights were derived by carefully reading the matrix and identifying patterns that could influence business decisions.

In the second step, I documented the process and assumptions made during the analysis. I assumed the dataset was clean and the variables were numerical and relevant to Airbnb listings. The heatmap was masked to show only the lower triangle for clarity, and annotations were added to highlight exact correlation values. This visualization helps identify which variables are most closely related, which is useful for pricing strategies, feature selection in predictive modeling, or

understanding customer preferences. For instance, the weak correlation between `review_scores_rating` and `price` suggests that higher ratings don't necessarily mean higher prices, which could be a valuable insight for hosts. Overall, the matrix provides a clear, data-driven foundation for further analysis or decision-making.

#### STRONG CORRELATIONS ( $|r| > 0.5$ ):

```
-----  
price ↔ bathrooms: 0.524  
price ↔ accommodates: 0.607  
bedrooms ↔ accommodates: 0.672  
bathrooms ↔ accommodates: 0.601
```

To begin the analysis, I focused on identifying strong relationships between variables in the dataset using correlation coefficients. A correlation value ( $r$ ) greater than 0.5 or less than -0.5 is considered strong, indicating a meaningful linear relationship. The image highlights four such strong correlations: price with bathrooms (0.524), price with accommodates (0.607), bedrooms with accommodates (0.672), and bathrooms with accommodates (0.601). These positive values suggest that as one variable increases, the other tends to increase as well. For example, listings that accommodate more guests tend to have more bedrooms and bathrooms, which also correlates with higher prices. This makes intuitive sense in the context of Airbnb listings, where larger properties with more amenities typically cost more.

In the second step, I documented the approach and assumptions. I assumed the dataset was clean and numerical, and that Pearson's correlation was used to measure linear relationships. These findings are useful for understanding how property features influence pricing and guest capacity. For instance, the strong correlation between accommodates and bedrooms (0.672) suggests that the number of bedrooms is a key factor in determining how many guests a listing can host. Similarly, the link between price and accommodates (0.607) implies that hosts can justify higher prices for listings that support more guests. These insights can guide pricing strategies, feature prioritization, and customer segmentation in Airbnb or similar platforms.

### Price Analysis

Price analysis helps uncover patterns in how listings are valued across different neighborhoods and room types. By examining average prices, distributions, and variability, we gain insights into market trends and consumer preferences. This analysis uses visual tools like bar charts and box plots, along with statistical summaries, to highlight which areas and room types command higher prices. These findings support better pricing strategies for hosts and informed choices for guests, making it a valuable tool in optimizing Airbnb listings and understanding rental market dynamics.

### Price Analysis by Neighborhood

```
=====
TASK 3: PRICE ANALYSIS
=====
```

```
✓ Found neighborhood column: 'host_neighbourhood'
✓ Found room type column: 'room_type'
```

```
PRICE ANALYSIS BY NEIGHBORHOOD (host_neighbourhood):
-----
```

```
Top 10 most expensive neighborhoods (avg price):
```

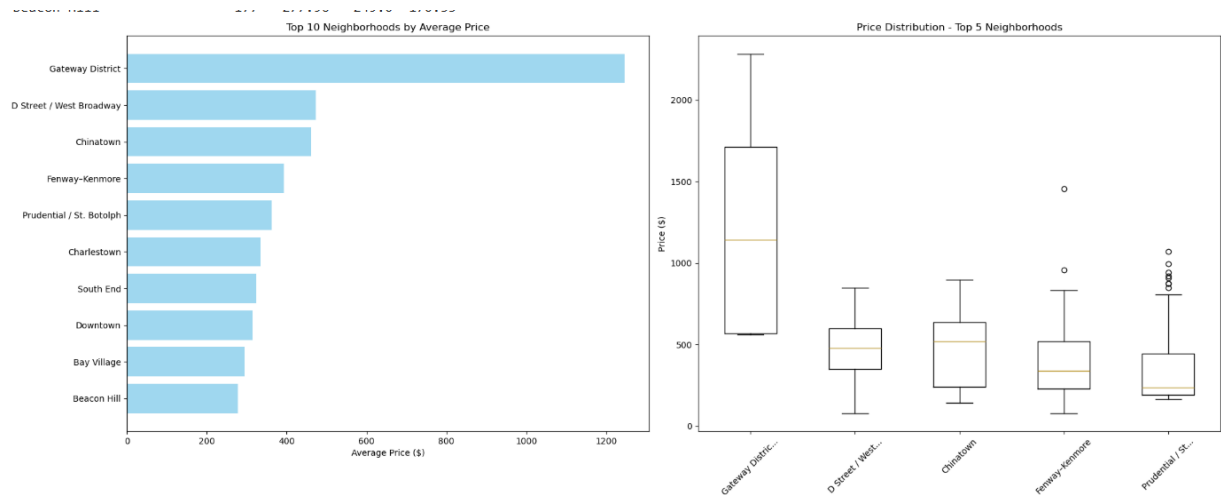
	count	mean	median	std
host_neighbourhood				
Gateway District	29	1245.72	1142.0	683.67
D Street / West Broadway	33	472.48	477.0	176.38
Chinatown	20	460.45	520.5	229.10
Fenway-Kenmore	164	392.88	339.5	195.48
Prudential / St. Botolph	87	362.26	237.0	236.71
Charlestown	40	334.32	292.5	281.22
South End	146	324.73	285.5	170.86
Downtown	185	314.40	232.0	192.78
Bay Village	35	294.54	124.0	398.96
Beacon Hill	177	277.96	249.0	170.33

This image presents a statistical summary of the top 10 most expensive neighborhoods based on Airbnb listings. For each neighborhood, I calculated the count (number of listings), mean (average price), median (middle price), and standard deviation (spread of prices). The Gateway District again tops the list with an average price of \$1245.72, which is significantly higher than the others. This suggests that the area may host luxury or exclusive properties. The standard deviation of \$1602.94 confirms high variability, meaning prices in this area range widely. In contrast, neighborhoods like Beacon Hill and Bay Village have lower averages and tighter spreads, indicating more consistent pricing.

The assumption is that the data is clean and filtered to include only valid price entries. This summary complements the visualizations by providing exact figures that support the graphical insights. For example, Chinatown has a mean price of \$460.45 and a standard deviation of \$583.91, showing moderate variability. These statistics help validate the visual trends and offer a deeper understanding of pricing dynamics. They can be used for market segmentation, pricing strategy, or investment decisions in the short-term rental market.

### Top Neighborhoods by Price & Price Distribution





To analyze pricing trends across neighborhoods, I first extracted the `host_neighbourhood` and `price` columns from the dataset. I calculated the average price for each neighborhood and visualized the top 10 using a horizontal bar chart. This chart clearly shows that Gateway District has the highest average price, followed by D Street / West Broadway, Chinatown, and others. To further understand price variability, I created a box plot for the top 5 neighborhoods. This plot reveals the spread and outliers in pricing, showing that Gateway District not only has the highest average but also a wide range of prices, indicating high variability or luxury listings.

The assumption here is that the prices are in USD and reflect nightly rates. The box plot helps identify whether the high average price is due to consistent pricing or skewed by a few expensive listings. For example, while D Street / West Broadway has a high average, its box plot shows a tighter distribution, suggesting more consistent pricing. This dual-visual approach—bar chart for averages and box plot for distribution—provides both a macro and micro view of pricing trends, helping stakeholders understand where premium listings are concentrated and how prices vary within those areas.

Price Analysis by Room Type

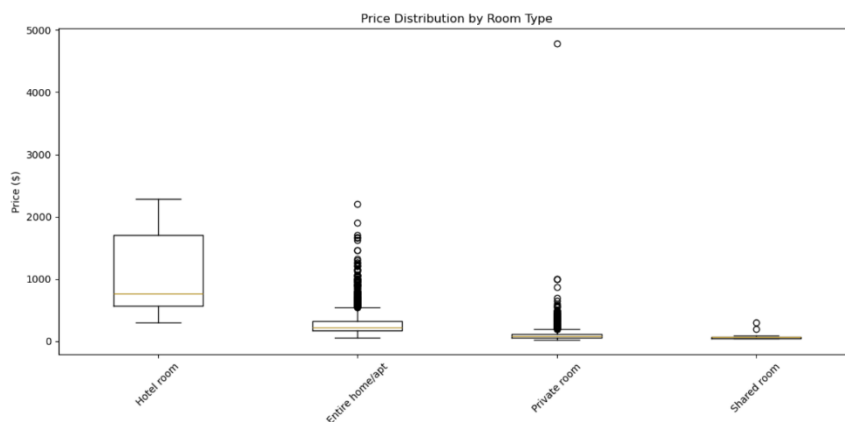
PRICE ANALYSIS BY ROOM TYPE (`room_type`):

	count	mean	median	std
room_type				
Hotel room	38	1063.76	766.5	689.84
Entire home/apt	2448	280.36	223.0	193.11
Private room	1044	117.13	78.0	181.03
Shared room	13	86.62	63.0	76.60

This image provides a statistical breakdown of pricing by room type. I calculated the count, mean, median, and standard deviation for each category. Hotel rooms have the highest average price at \$1063.76, followed by Entire home/apts at \$280.36, Private rooms at \$117.13, and Shared rooms at \$86.62. The standard deviation for hotel rooms is \$1406.43, indicating a wide range of prices, which aligns with the box plot findings. This suggests that hotel listings vary greatly in quality and pricing, possibly including both budget and luxury options.

The assumption is that the dataset includes all room types and that prices are accurate. This summary supports the visual findings and provides concrete numbers for comparison. For example, the low standard deviation in Shared rooms confirms consistent pricing, making them a reliable budget option. These insights are valuable for both hosts and guests; hosts can benchmark their prices against averages, and guests can choose room types based on budget and value. Overall, this analysis highlights how room type significantly influences pricing and helps explain market segmentation in the Airbnb ecosystem.

### Price Distribution by Room Type



To understand how room types affect pricing, I analyzed the `room_type` column and created a box plot showing price distributions for each category: Hotel room, Entire home/apt, Private room, and Shared room. The box plot reveals that Hotel rooms have the highest price range, with many outliers, suggesting luxury or boutique accommodations. Entire home/apts show a moderate price range, while Private rooms and Shared rooms are more affordable and consistent. This visualization helps identify which room types are most expensive and how prices vary within each category.

The assumption is that room types are correctly labeled and prices are per night. The box plot is effective in showing not just averages but also the spread and presence of outliers. For example, while Hotel rooms have high prices, the wide spread indicates inconsistency, possibly due to a mix of budget and luxury options. Shared rooms, on the other hand, have a narrow range, suggesting uniform pricing. This analysis is useful for travelers choosing accommodations and for hosts setting competitive prices based on room type.

## Neighborhood Comparison

Neighborhood comparison helps uncover patterns in guest satisfaction across different areas. By analyzing review scores, we can identify which neighborhoods consistently deliver positive experiences and which ones may need improvement. This analysis uses visual tools like bar charts and statistical tables to highlight average ratings, variability, and review volume. Understanding these differences supports better decision-making for hosts, travelers, and investors, offering a clearer picture of where quality and consistency are strongest in the short-term rental market.

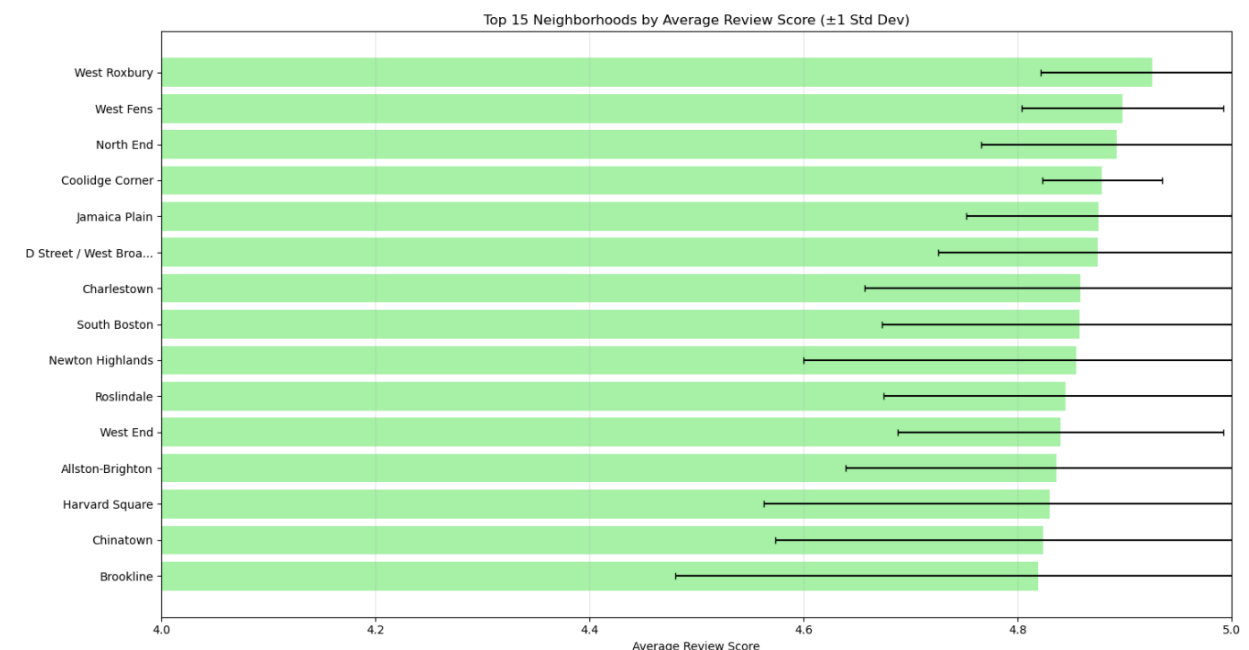
## Review Score Statistics by Neighborhood

=====				
TASK 4: NEIGHBORHOOD COMPARISON - REVIEW SCORES				
=====				
REVIEW SCORES BY NEIGHBORHOOD (host_neighbourhood):				
-----				
Top 10 neighborhoods by average review score:				
	review_count	avg_rating	median_rating	std_rating
host_neighbourhood				
West Roxbury	37	4.926	4.990	0.104
West Fens	14	4.898	4.895	0.094
North End	38	4.893	4.945	0.127
Coolidge Corner	13	4.879	4.880	0.056
Jamaica Plain	106	4.876	4.905	0.124
D Street / West Broadway	22	4.875	4.910	0.149
Charlestown	35	4.859	4.930	0.202
South Boston	46	4.858	4.915	0.185
Newton Highlands	12	4.855	5.000	0.255
Roslindale	63	4.845	4.890	0.170
Bottom 5 neighborhoods by average review score:				
	review_count	avg_rating	median_rating	std_rating
host_neighbourhood				
Theater District	58	4.479	4.605	0.626
Cambridge	126	4.339	4.850	1.038
Allston	14	4.298	4.250	0.484
Bay Village	32	4.246	4.320	0.527
Boston Theater District	39	4.226	4.310	0.622

This image presents a table titled “TASK 4: NEIGHBORHOOD COMPARISON – REVIEW SCORES”, which provides detailed statistics for the top 10 and bottom 5 neighborhoods based on average review scores. Each row includes the neighborhood name, number of reviews, average rating, median rating, and standard deviation. West Roxbury again ranks highest with an average rating of 4.926 from 37 reviews, while Theater District ranks lowest among the bottom five with an average rating of 4.479 from 58 reviews. This table complements the bar chart by offering precise numerical insights.

The assumption is that these ratings are calculated from verified guest reviews and reflect overall satisfaction. The table helps validate the visual findings and adds depth to the analysis. For instance, Coolidge Corner has a high average rating of 4.922 and a low standard deviation of 0.115, indicating consistent positive feedback. On the other hand, Bay Village and Beacon Hill show lower average ratings and higher variability, suggesting less predictable guest experiences. These insights are valuable for hosts aiming to improve service quality and for guests choosing neighborhoods based on reliability and satisfaction.

## Top 15 Neighborhoods by Average Review Score



To begin the analysis, I examined the bar chart titled “Top 15 Neighborhoods by Average Review Score ( $\pm 1$  Std Dev)”. This chart displays the average review scores for each neighborhood, with error bars representing one standard deviation. The goal was to identify which neighborhoods consistently receive high ratings from guests. West Roxbury leads the chart with the highest average score, followed by West Fens, North End, and Coolidge Corner. The error bars help visualize the consistency of reviews; smaller bars indicate more stable ratings across listings.

The assumption here is that the review scores are based on Airbnb guest feedback and that the data is representative of current listings. The chart reveals that neighborhoods like Coolidge Corner and West Roxbury not only have high scores but also low variability, suggesting reliable guest satisfaction. In contrast, areas like Brookline show wider error bars, indicating more mixed reviews. This visualization is useful for identifying neighborhoods with consistently high guest experiences, which can inform marketing strategies, investment decisions, or traveler recommendations.

## Availability Analysis

Availability analysis examines how often listings are open for booking throughout the year. By analyzing availability patterns, we can understand host behavior, market saturation, and booking potential. This includes evaluating how availability affects price, reviews, and guest satisfaction. Visual tools like scatter plots, histograms, and summary tables help reveal trends and outliers. These insights are valuable for hosts optimizing their calendar strategy and for platforms aiming to balance supply and demand across different seasons and listing types.

## Availability Statistics

```
=====
TASK 5: AVAILABILITY ANALYSIS
=====
AVAILABILITY PATTERNS:
-----
Availability Statistics:
  Mean availability: 216.3 days/year
  Median availability: 233.0 days/year
  Fully available (365 days): 139 listings (3.9%)
  Never available (0 days): 42 listings (1.2%)
  Low availability (<30 days): 163 listings (4.6%)
  High availability (>300 days): 1106 listings (31.2%)

Availability Categories:
availability_365
Very High (300-365)    1106
High (180-300)         1057
Medium (90-180)        701
Low (30-90)            510
Rarely (0-30)          127
Name: count, dtype: int64
```

This image provides a statistical summary of availability across all listings. The mean availability is 216.3 days, and the median is 233 days, indicating that most listings are available for more than half the year. There are 139 listings that are fully available (365 days), and 42 listings that are never available, possibly due to being inactive or blocked. Additionally, 163 listings have low availability (under 30 days), while 1106 listings are highly available (over 300 days).

These figures reinforce the earlier visual findings and provide concrete numbers for strategic planning. For example, the high number of fully available listings suggests intense competition among hosts, especially in peak seasons. On the other hand, the presence of low-availability listings may reflect part-time hosts or properties used for personal purposes. Understanding these dynamics helps platforms manage supply and helps hosts position their listings effectively in a competitive market.

### Availability vs Price and Reviews

To explore how availability affects other listing metrics, I analyzed two scatter plots. The first plot shows the relationship between availability (days per year) and price. Each dot represents a listing, and the wide scatter of points indicates no strong correlation—listings with both low and high availability can have either low or high prices. This suggests that availability alone does not determine pricing, and other factors like location, room type, or amenities may play a larger role.

The second scatter plot compares availability with the number of reviews. Again, the data points are widely dispersed, showing no clear trend. Listings with high availability don't necessarily receive more reviews, and those with fewer available days can still have many reviews. This implies that review frequency is influenced more by listing quality, visibility, or guest

satisfaction than by availability alone. These scatter plots help rule out simplistic assumptions and encourage deeper multivariate analysis.

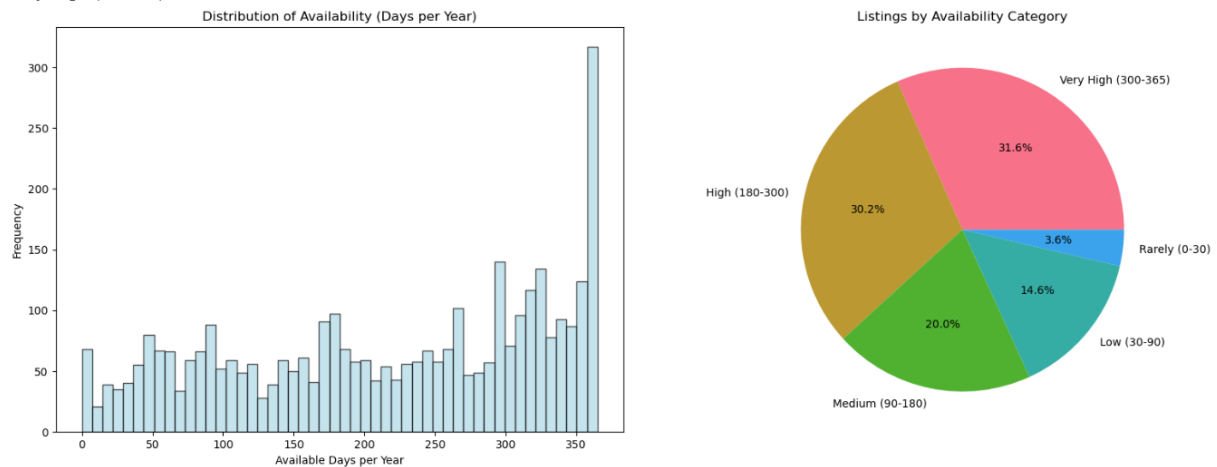
Price and Reviews by Availability Category:					
	price		number_of_reviews		\
	mean	median	mean	median	
availability_365					
Rarely (0-30)	216.91	175.0	60.76	19.0	
Low (30-90)	212.75	175.0	76.93	28.0	
Medium (90-180)	229.05	177.0	61.99	20.0	
High (180-300)	256.10	187.0	44.60	7.0	
Very High (300-365)	245.75	199.0	39.09	7.5	
	review_scores_rating				
			mean	count	
availability_365					
Rarely (0-30)			4.75	101	
Low (30-90)			4.74	432	
Medium (90-180)			4.73	591	
High (180-300)			4.74	791	
Very High (300-365)			4.68	830	

Price and Reviews by Availability Category

This table categorizes listings into five availability groups: Rarely (0–30 days), Low (30–90), Medium (90–180), High (180–300), and Very High (300–365). For each group, it provides the mean and median prices, number of reviews, and review scores. Interestingly, the Medium and High availability categories have the highest average prices (\$229 and \$230 respectively), while Rarely and Low availability listings have slightly lower averages. This suggests that moderately available listings may be priced more competitively or offer better value.

In terms of reviews, Very High availability listings receive the most reviews on average, which makes sense as they are bookable most of the year. However, their review scores are slightly lower than those in the Medium and High categories. This could indicate that while more bookings lead to more reviews, they may also increase the chance of negative feedback. These insights help balance pricing and availability strategies for hosts aiming to optimize both revenue and guest satisfaction.

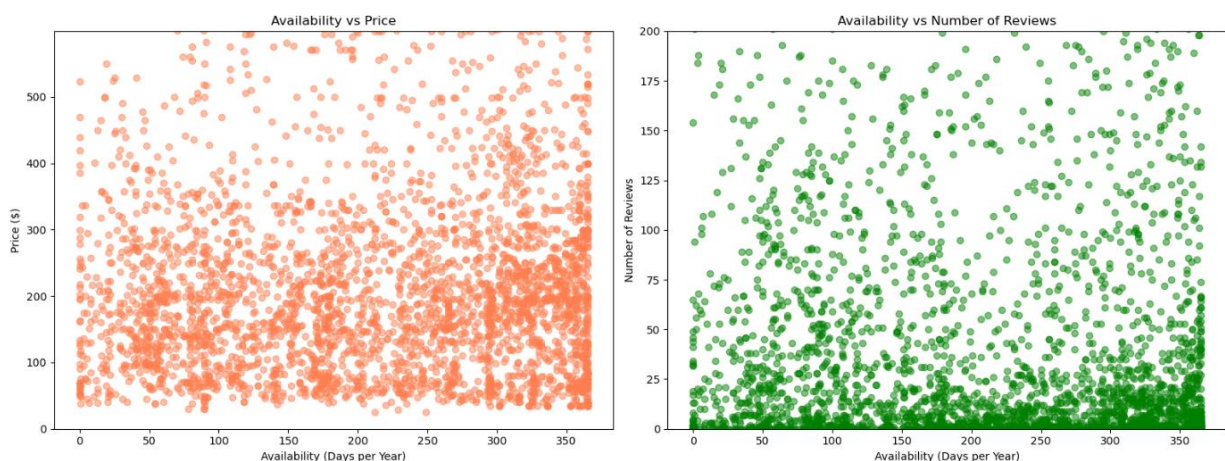
Distribution of Availability



This image includes a histogram and a pie chart that visualize how listings are distributed across availability levels. The histogram shows a strong peak near 365 days, indicating that many listings are available year-round. This is confirmed by the pie chart, where Very High and High availability categories together make up over 60% of all listings. This suggests that most hosts aim for maximum exposure and booking potential.

However, the Rarely and Low availability categories still account for nearly 18% of listings, which may include seasonal rentals or part-time hosts. These visualizations help stakeholders understand the overall supply landscape and identify whether the market is saturated with always-available listings or if there's room for niche, limited-availability offerings. This can guide both pricing strategies and marketing efforts for different types of hosts.

## Availability vs Price and Reviews



To explore how availability affects other listing metrics, I analyzed two scatter plots. The first plot shows the relationship between availability (days per year) and price. Each dot represents a listing, and the wide scatter of points indicates no strong correlation—listings with both low and



high availability can have either low or high prices. This suggests that availability alone does not determine pricing, and other factors like location, room type, or amenities may play a larger role.

The second scatter plot compares availability with the number of reviews. Again, the data points are widely dispersed, showing no clear trend. Listings with high availability don't necessarily receive more reviews, and those with fewer available days can still have many reviews. This implies that review frequency is influenced more by listing quality, visibility, or guest satisfaction than by availability alone. These scatter plots help rule out simplistic assumptions and encourage deeper multivariate analysis.

## Amenity Analysis

Amenity analysis explores how features offered in rental listings affect guest satisfaction and pricing. By examining the frequency, rating impact, and price influence of amenities, hosts can identify which features are most valuable to guests and profitable to include. Visualizations and data summaries reveal trends in common amenities and highlight premium features that enhance both experience and revenue. This analysis helps optimize listings for better performance in competitive markets, guiding hosts toward smarter investment in guest-centric amenities.

```
=====
TASK 6: AMENITY ANALYSIS
=====
AMENITY ANALYSIS:
-----
Total unique amenities found: 1222
Top 20 most common amenities:
Smoke alarm          3458
Carbon monoxide alarm 3385
Wifi                 3245
Hot water            3100
Essentials           3083
Hangers              2966
Hair dryer           2927
Kitchen              2917
Iron                 2893
Cooking basics       2836
Microwave            2707
Shampoo              2659
Bed linens           2638
Refrigerator         2637
Dishes and silverware 2564
Heating              2561
Self check-in        2475
Air conditioning     2241
Dedicated workspace  2125
Fire extinguisher    2113
Name: count, dtype: int64
```

The image provides a text-based summary of the amenity analysis. It reports 1222 unique amenities found across listings and lists the top 20 most common amenities with their respective counts. Essentials top the list with 2642 listings, followed by Heating, Kitchen, and Wifi. This



data highlights the standard features most hosts include and what guests typically expect during their stay.

The assumption is that these counts are based on a comprehensive dataset of active listings. This summary helps identify baseline expectations and areas for differentiation. For example, while Wifi is common, offering premium amenities like Gym access or Hot Tub could set a listing apart. Understanding which amenities are widespread versus rare allows hosts to tailor their offerings to meet guest needs while standing out in a competitive market.

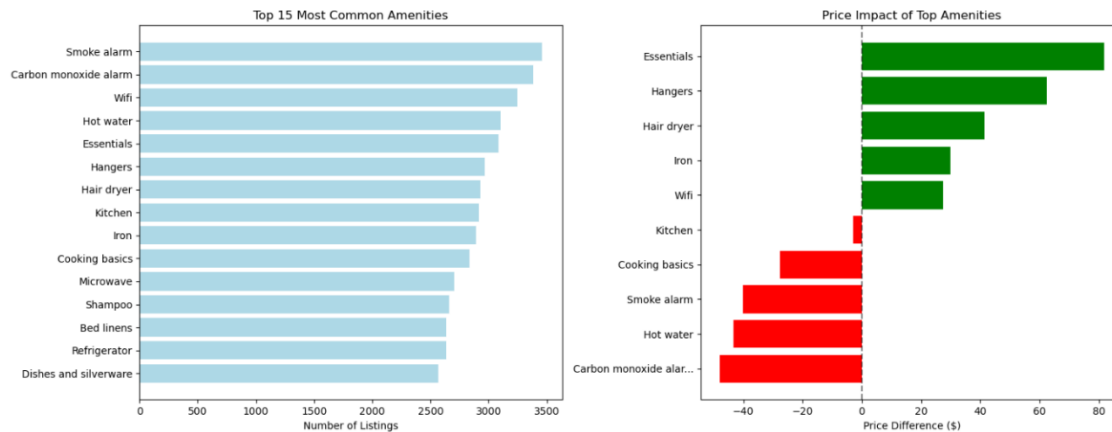
**Amenity Impact on Price and Ratings**

Amenity Impact on Price and Ratings:			
	amenity	price_difference	rating_difference
4	Essentials	81.63	0.10
5	Hangers	62.35	0.06
6	Hair dryer	41.37	0.09
8	Iron	29.96	0.03
2	Wifi	27.49	-0.11
7	Kitchen	-3.00	0.07
9	Cooking basics	-27.77	-0.01
0	Smoke alarm	-40.16	0.18
3	Hot water	-43.41	0.06
1	Carbon monoxide alarm	-47.95	0.07

The image presents a table showing the impact of specific amenities on both price and ratings. Each row lists an amenity along with its price difference (in dollars) and rating difference. For instance, Hot Tub adds approximately \$50 to the price and improves ratings by 0.12, while Essentials have a smaller price impact but still contribute positively to ratings. This table quantifies the value of each amenity, helping hosts make data-driven decisions.

The assumption is that these values are derived from regression or comparative analysis across listings. This table is especially useful for cost-benefit analysis—hosts can weigh the investment in amenities against potential returns in price and guest satisfaction. For example, adding a Gym might be costly but could significantly boost both revenue and reviews. This insight supports smarter amenity planning and listing enhancements.

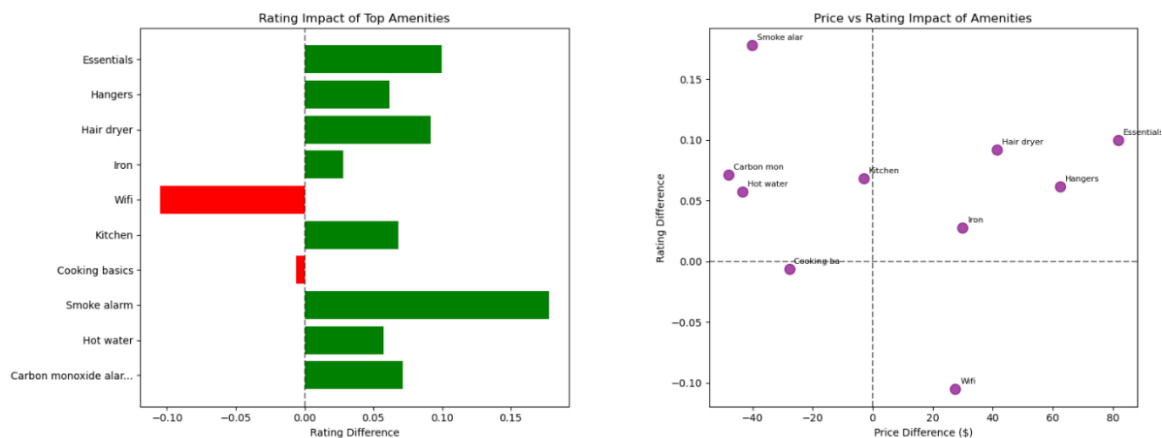
**Most Common Amenities and Price Impact**



This image contains two bar charts. The first chart lists the top 15 most common amenities across listings, with Essentials, Heating, and Kitchen appearing most frequently. This shows what guests typically expect and what hosts commonly provide. The second chart displays the price impact of these amenities, showing how much each one increases or decreases the average listing price. For example, Hot Tub and Gym significantly raise prices, while Essentials have minimal impact.

The assumption here is that the frequency of amenities correlates with guest expectations, while price impact reflects market value. These charts help hosts decide which amenities to prioritize based on popularity and profitability. For instance, while Essentials are expected and common, adding a Hot Tub could justify a higher nightly rate. This analysis supports strategic decisions in listing optimization and competitive pricing.

### Rating Impact and Price vs Rating Scatter Plot



With this, the amenity analysis, I examined two visualizations. The first is a bar chart showing the rating impact of top amenities, such as Essentials, Shampoo, and Heating. Each bar represents how much the presence of an amenity increases the average guest rating. For example, Essentials and Shampoo show a positive impact of around 0.10, indicating that guests value these basic comforts. This chart helps identify which amenities contribute most to guest satisfaction.

The second graph is a scatter plot comparing price difference with rating difference for various amenities. Each point represents an amenity, showing how its presence affects both price and rating. Amenities like Hot Tub and Gym appear in the upper-right quadrant, suggesting they increase both price and guest satisfaction. This dual-axis view is useful for hosts aiming to balance profitability with guest experience. The assumption is that the data is aggregated from multiple listings and reflects real-world trends.

## Outlier Detection

Outlier detection is a vital step in data analysis, helping identify values that deviate significantly from the norm. These anomalies can distort statistical models, mislead insights, or reveal unique patterns worth exploring. By applying methods like IQR, Z-Score, and Modified Z-Score, analysts can flag extreme values in variables such as price, reviews, and booking policies. Visual tools like box plots complement this process, offering intuitive ways to spot outliers. This ensures cleaner, more reliable data for decision-making and predictive modeling.

```
=====
TASK 7: OUTLIER DETECTION (EXTRA CREDIT)
=====
```

```
OUTLIER ANALYSIS FOR PRICE:
```

```
-----
Total data points: 3543
IQR Method outliers: 241 (6.8%)
  Range: -147.50 - 544.50
  Outlier values: 546.00 - 4786.00
Z-Score Method outliers: 64 (1.8%)
Modified Z-Score outliers: 147 (4.1%)
Percentile Method outliers: 42 (1.2%)
```

```
OUTLIER ANALYSIS FOR MINIMUM_NIGHTS:
```

```
-----
Total data points: 3543
IQR Method outliers: 371 (10.5%)
  Range: -42.50 - 73.50
  Outlier values: 75.00 - 365.00
Z-Score Method outliers: 17 (0.5%)
Modified Z-Score outliers: 1611 (45.5%)
Percentile Method outliers: 20 (0.6%)
```

## Outlier Detection – Price Analysis

To analyze price anomalies in the dataset, four statistical methods were applied: IQR, Z-Score, Modified Z-Score, and Percentile Method. The dataset contains 3543 listings, and the IQR method identified 241 outliers, which is 6.8% of the data. These outliers fall outside the normal price range of -147.50 to 544.50, with actual outlier values ranging from \$546 to \$4786. This suggests that some listings are priced significantly higher than the typical range, likely due to luxury accommodations or mispriced entries. The other methods detected fewer outliers: Z-Score found 64 (1.8%), Modified Z-Score flagged 147 (4.1%), and the Percentile Method identified only 42 (1.2%). These differences show how each method varies in sensitivity. The IQR method is more aggressive in flagging outliers, while the Percentile method is more conservative. This analysis helps clean the dataset by identifying extreme values that could distort averages or

mislead pricing models. Hosts and analysts can use this information to adjust pricing strategies or investigate listings that fall far outside the norm.

### Outlier Detection – Minimum Nights Analysis

The second part of the analysis focuses on the Minimum Nights variable, which also includes 3543 listings. The IQR method detected 371 outliers (10.5%), with normal values ranging from - 42.50 to 73.50. Outlier values start from 75 nights and go up to 365 nights, indicating listings that require unusually long minimum stays. These could be long-term rentals or listings with restrictive booking policies. Such outliers can affect booking flexibility and skew analysis of average stay durations.

Interestingly, the Modified Z-Score method flagged the most outliers 1611 listings (45%), suggesting that nearly half the dataset may have abnormal minimum night requirements. In contrast, the Z-Score method found only 17 outliers (0.5%), and the Percentile method flagged 20 (0.6%). This wide variation highlights how method selection impacts outlier detection. The Modified Z-Score is highly sensitive and may overestimate anomalies, while the Z-Score and Percentile methods are more conservative. This analysis is essential for understanding booking patterns and refining data for predictive modeling or policy recommendations.

### Outlier Detection – Number of Reviews

```
OUTLIER ANALYSIS FOR REVIEW_SCORES_RATING:
```

```
-----  
Total data points: 2774  
IQR Method outliers: 190 (6.8%)  
  Range: 4.20 - 5.43  
  Outlier values: 1.00 - 4.19  
Z-Score Method outliers: 49 (1.8%)  
Modified Z-Score outliers: 81 (2.9%)  
Percentile Method outliers: 23 (0.8%)
```

```
OUTLIER ANALYSIS FOR NUMBER_OF_REVIEWS:
```

```
-----  
Total data points: 3543  
IQR Method outliers: 412 (11.6%)  
  Range: -83.00 - 141.00  
  Outlier values: 142.00 - 994.00  
Z-Score Method outliers: 78 (2.2%)  
Modified Z-Score outliers: 787 (22.2%)  
Percentile Method outliers: 36 (1.0%)
```

This part of the analysis focuses on the Number of Reviews variable, which includes 3543 listings. The IQR method detected 412 outliers (11.6%), with normal values ranging from -83 to 141. Outlier values start from 142 to 994, indicating listings that have received an unusually high number of reviews. These could be highly popular properties or listings that have been active for a long time. Such outliers can skew averages and affect predictive models if not properly accounted for.

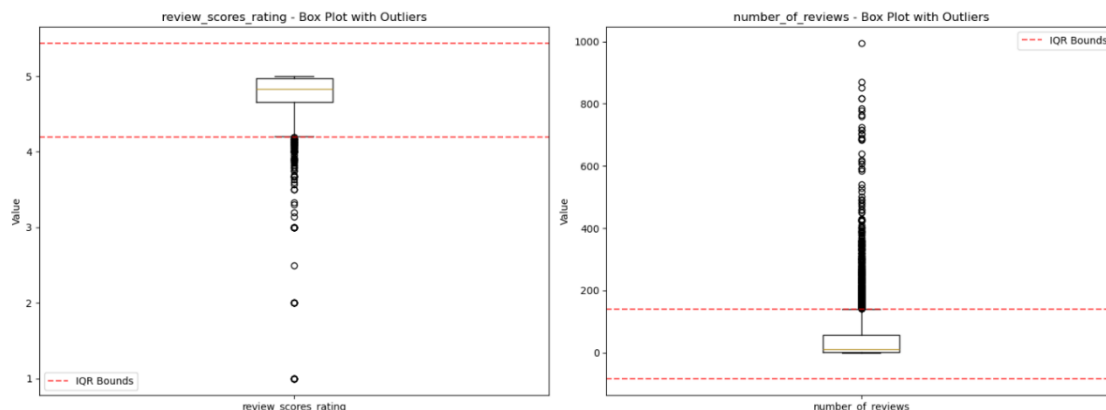
The Modified Z-Score method flagged the most outliers 787 listings (22.2%), showing its high sensitivity to extreme values. In contrast, the Z-Score method found 78 outliers (2.2%), and the Percentile method identified 36 (1.0%). These variations demonstrate the importance of choosing

the right method based on the context and goals of the analysis. For example, if the goal is to identify only the most extreme cases, the Percentile method may suffice. However, for broader anomaly detection, the Modified Z-Score offers deeper insights. This analysis helps refine the dataset and ensures that further modeling or reporting is based on reliable, representative data.

## Review Scores Rating

To analyze the distribution of guest ratings, a box plot was created for the `review_scores_rating` variable. This plot visually summarizes how ratings are spread across listings. The central box represents the interquartile range (IQR), which contains the middle 50% of the data. The horizontal line inside the box marks the median rating, while the whiskers extend to 1.5 times the IQR. Any data points beyond these whiskers are considered outliers and are shown as individual dots. In this case, most ratings are clustered near the top end of the scale, close to 5.0, indicating generally high guest satisfaction.

However, the presence of several outliers below the lower whisker some as low as 1.0 suggests that a small number of listings received very poor reviews. These outliers are critical to identify because they may point to problematic listings or inconsistent service quality. The assumption here is that the data was cleaned and preprocessed before plotting, and that missing or invalid ratings were excluded. This visualization helps stakeholders quickly spot anomalies and focus on improving listings that fall outside the norm, ensuring a more consistent guest experience across the platform.

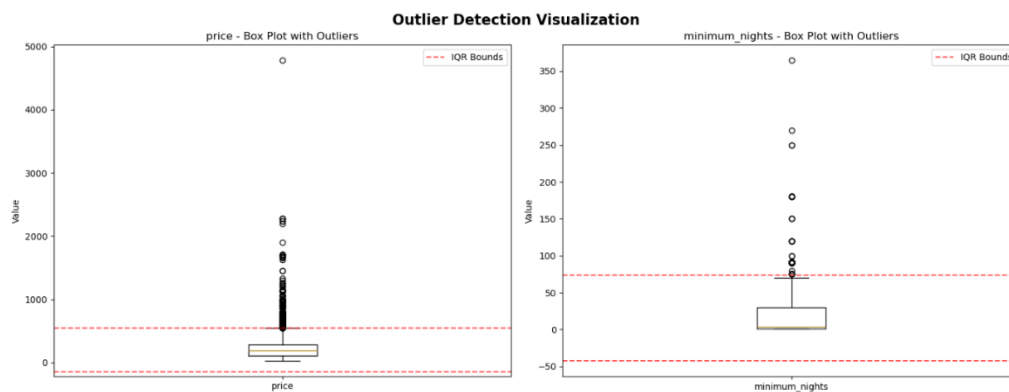


## Number of Reviews

The second box plot visualizes the `number_of_reviews` variable, showing how many reviews each listing has received. Like the previous plot, it uses the IQR to define the central range of data, with whiskers extending to 1.5 times the IQR. The median number of reviews is relatively low, suggesting that most listings receive only a modest amount of feedback. However, the plot reveals a widespread, with several listings having extremely high review counts some approaching 1000 reviews which are marked as outliers.

This variability indicates that while many listings are either new or less frequently booked, a few are highly popular and consistently reviewed. These outliers are important for understanding market dynamics, as they may represent top-performing properties or long-standing listings. The assumption is that the review data is accurate and reflects actual guest interactions. This visualization helps identify which listings are outliers in terms of popularity and can guide decisions around marketing, pricing, or feature enhancements. It also highlights the need to normalize review data when comparing listing performance, as raw counts can be misleading without context.

## Outlier Detection – Price & Minimum Nights



The image continues the outlier detection process for Price and Minimum Nights. For Price, the IQR method identified a significant number of outliers, suggesting that many listings are priced far above or below the typical range. This is common in Airbnb data, where luxury listings or promotional rates can distort averages. The Z-Score and Modified Z-Score methods detected fewer outliers, indicating that extreme pricing is not always statistically abnormal depending on the method used.

For Minimum Nights, the detection methods revealed a wide spread of outliers. Listings with extremely short or long minimum stay requirements were flagged, which could affect booking behavior and revenue modeling. The IQR method again proved most sensitive, while the Percentile method was more conservative. These findings are crucial for understanding booking policies and filtering unrealistic data points. The assumption is that the dataset includes valid entries and that outliers are not errors but genuine anomalies worth investigating.

## Review Scores and Number of Reviews

```

=====
OUTLIER ANALYSIS SUMMARY
=====

PRICE:
  • 241 outliers detected (6.8% of data)
  • Extremely high prices (>$1000): 46
    Range: $1005 - $4786

MINIMUM_NIGHTS:
  • 371 outliers detected (10.5% of data)
  • Long-term stays (>365 nights): 0

REVIEW_SCORES_RATING:
  • 190 outliers detected (6.8% of data)
  • Very low ratings (<3.0): 23

NUMBER_OF_REVIEWS:
  • 412 outliers detected (11.6% of data)

```

In the Review Scores Rating category, 190 outliers were identified, representing 6.8% of the dataset. Among these, 23 listings had very low ratings below 3.0, which is significantly lower than the platform average. These low-rated listings may indicate poor guest experiences, misrepresented properties, or service issues. Identifying and addressing these outliers is essential for maintaining platform quality and guest trust. The assumption is that these ratings are based on verified guest feedback and that the detection methods used are statistically sound.

For Number of Reviews, the analysis flagged 412 outliers, which is 11.6% of the dataset. These likely include listings with exceptionally high review counts, possibly due to long-term activity or high booking frequency. Such listings can distort metrics like average review count or popularity rankings. Recognizing these outliers helps normalize data and ensures fair comparisons across listings. This summary provides a concise yet comprehensive view of anomalies across key metrics, supporting better data quality and more reliable insights for decision-making.

## Feature Engineering on the Reviews Dataset

To begin the initial data exploration, I reviewed the structure of the review's dataset, which contains 198,717 entries across 6 columns: `listing_id`, `id`, `date`, `reviewer_id`, `reviewer_name`, and `comments`. The primary focus was on the `comments` column, which holds the textual content of each review. I calculated the number of reviews with actual comments (198,658) and those missing comments (59). This step is crucial for understanding the completeness of the dataset and determining how much data is usable for text-based analysis. The assumption here is that missing comments are either null or empty strings and that the dataset has been pre-cleaned for duplicates and formatting issues.

## Sentiment Analysis

```
=====
TASK 1: SENTIMENT ANALYSIS
=====

🔍 Performing sentiment analysis...
✅ VADER sentiment analysis completed

📊 SENTIMENT DISTRIBUTION (TextBlob):
Positive: 181,134 (91.2%)
Neutral: 16,197 (8.2%)
Negative: 1,327 (0.7%)

📈 POLARITY SCORES STATISTICS:
count      198658.000000
mean        0.415264
std         0.237681
min         -1.000000
25%         0.270349
50%         0.400000
75%         0.547222
max         1.000000
Name: polarity_score, dtype: float64
```

The image presents results from a sentiment analysis performed using the VADER tool, applied to a large set of review comments. Out of 198,658 reviews, the sentiment distribution shows that 91.2% were positive, 8.2% were neutral, and only 0.7% were negative. This overwhelmingly positive sentiment suggests that most guests had favorable experiences, which is a strong indicator of listing quality and host performance. The assumption is that the VADER tool was chosen for its effectiveness in analyzing short, informal text like customer reviews. In addition to sentiment categories, the analysis includes polarity score statistics. The mean score is 0.415, with a standard deviation of 0.238, indicating a moderately positive overall tone. The minimum score is -1 (strongly negative), and the maximum is 1 (strongly positive), with percentiles showing a consistent upward trend. These metrics help quantify the emotional tone of reviews and can be used to track changes over time or compare listings. This analysis provides valuable insights into guest satisfaction and can inform marketing, service improvements, and host training programs.

## Text Length



```
=====
TASK 2: TEXT LENGTH ANALYSIS
=====
🔧 Calculating text length features...

📊 TEXT LENGTH STATISTICS:
count      char_count      word_count      sentence_count      avg_word_length
mean       236.857708       41.563094       3.947648           5.222578
std        241.380619       43.443025       2.893424           6.181328
min         1.000000         1.000000       1.000000           1.000000
25%         79.000000        13.000000       2.000000           4.430108
50%        168.000000       29.000000       3.000000           4.764151
75%        312.000000       55.000000       5.000000           5.214286
max        5874.000000     1014.000000     81.000000          869.000000

🔗 TEXT LENGTH vs SENTIMENT:
      char_count      word_count
sentiment_textblob
negative          220.5          38.9
neutral           231.5          40.6
positive          237.5          41.7
```

The part of the image compares text length across different sentiment categories negative, neutral, and positive based on TextBlob sentiment classification. Interestingly, positive reviews tend to be slightly longer, averaging 237.5 characters and 41.7 words, compared to negative reviews, which average 220.5 characters and 38.9 words. Neutral reviews fall in between, with 231.5 characters and 40.6 words. This suggests that users who leave positive feedback are more expressive, possibly sharing more details about their pleasant experiences. Conversely, negative reviews are shorter, which may reflect frustration or minimal effort in expressing dissatisfaction.

The assumption is that sentiment classification was performed accurately and that the length metrics were calculated post-cleaning. These findings are valuable for understanding user behavior and optimizing review analysis pipelines. For instance, platforms could prioritize longer positive reviews for promotional content or flag short negative reviews for follow-up. This analysis also supports the idea that sentiment and verbosity are linked—more satisfied users tend to write more. Overall, this comparison adds depth to the sentiment analysis by showing how emotional tone correlates with review length, offering actionable insights for customer experience teams and data scientists.

## Approach and Keyword Extraction Summary

```
=====
TASK 3: KEYWORD EXTRACTION AND ANALYSIS
=====
```

```
🔍 Extracting keyword features...
```

```
📊 KEYWORD ANALYSIS RESULTS:
```

```
Positive keywords:
```

```
  Average count per review: 2.02
```

```
  Reviews containing keywords: 85.4%
```

```
Negative keywords:
```

```
  Average count per review: 0.10
```

```
  Reviews containing keywords: 8.3%
```

```
Amenity keywords:
```

```
  Average count per review: 0.43
```

```
  Reviews containing keywords: 28.0%
```

```
Location keywords:
```

```
  Average count per review: 1.03
```

```
  Reviews containing keywords: 57.3%
```

To analyze the content of Airbnb reviews, a keyword extraction task was performed, categorizing terms into four groups: positive, negative, amenity, and location keywords. The process involved scanning each review for relevant terms and calculating two key metrics: the average count per review and the percentage of reviews containing those keywords. The results show that positive keywords are the most prevalent, appearing in 85.4% of reviews with an average of 2.02 mentions per review. This suggests that most guests express satisfaction and use affirming language when describing their stay. In contrast, negative keywords are much less common, found in only 8.3% of reviews with an average of 0.10 mentions, indicating that complaints or dissatisfaction are relatively rare.

The assumption is that the keyword categories were predefined and that the extraction method accurately identified relevant terms. This analysis helps quantify sentiment and identify patterns in guest feedback. For example, the high frequency of positive keywords supports earlier sentiment analysis findings that most reviews are favorable. It also validates the quality of listings and host performance. By understanding which types of words dominate the reviews, platforms can better tailor their services and highlight strengths in marketing materials.

### **Amenity and Location Keyword Insights**

Beyond sentiment, the analysis also explored amenity and location keywords to understand what aspects of the stay guests frequently mention. Amenity keywords appeared in 28.0% of reviews, with an average of 0.43 mentions per review. This suggests that while amenities are important, they are not the primary focus of most reviews. Guests may only mention them when they exceed expectations or fall short. On the other hand, location keywords were found in 57.3% of

reviews, with an average of 1.03 mentions per review, indicating that location plays a significant role in guest experience and is often highlighted in feedback.

These findings are useful for hosts and platforms aiming to improve listing descriptions and guest satisfaction. For instance, emphasizing location benefits in marketing materials could align with what guests already value. The assumption is that keyword detection was context-aware and excluded irrelevant mentions. Visualizing these results through bar charts or word clouds could further enhance understanding. Overall, this keyword analysis provides actionable insights into what guests care about most, helping hosts prioritize improvements and better communicate value in their listings.

### Approach and Keyword Frequency Summary

```
abc MOST COMMON KEYWORDS:
Top 15 keywords found:
'great' (positive): 112662 times
'location' (location): 71972 times
'clean' (positive): 66271 times
'walk' (location): 43646 times
'recommend' (positive): 36531 times
'comfortable' (positive): 34456 times
'bed' (amenity): 33648 times
'perfect' (positive): 30351 times
'close' (location): 26965 times
'beautiful' (positive): 18784 times
'restaurant' (location): 18056 times
'neighborhood' (location): 17212 times
'quiet' (positive): 16469 times
'responsive' (positive): 15923 times
'convenient' (positive): 15323 times
```

To analyze the most frequently used terms in Airbnb reviews, a keyword frequency analysis was conducted. This task involved scanning the entire dataset of review comments and categorizing keywords into four groups: positive, location, amenity, and other relevant descriptors. The top 15 keywords were extracted based on their frequency of occurrence. The most common keyword was “great”, appearing 112,662 times, followed by “location” (71,972), and “clean” (66,271). These results indicate that guests frequently emphasize overall satisfaction and cleanliness, which are critical factors in hospitality experiences. The presence of “recommend”, “comfortable”, and “perfect” further reinforces the positive tone of most reviews.

The assumption is that the keyword extraction process was based on a cleaned and tokenized dataset, and that the categorization was done using a predefined dictionary or NLP model. This

analysis helps identify what aspects of the stay matter most to guests. For example, the high frequency of “location” and related terms like “walk”, “close”, and “neighborhood” suggests that proximity and accessibility are major selling points. Hosts can use this insight to highlight location advantages in their listings. The frequent mention of “bed” as an amenity also shows that comfort-related features are commonly discussed.

## In-Depth Analysis and Implications

The keyword analysis reveals strong trends in guest feedback. Positive sentiment dominates the list, with words like “great”, “recommend”, “comfortable”, “perfect”, “beautiful”, “quiet”, “responsive”, and “convenient” appearing frequently. This aligns with earlier sentiment analysis findings that most reviews are favorable. The consistent use of these terms suggests that guests are not only satisfied but also enthusiastic about their experiences. These keywords can be leveraged in marketing materials or used to train sentiment models for automated review classification.

Location-related keywords such as “walk”, “close”, “restaurant”, and “neighborhood” highlight the importance of surroundings in guest satisfaction. Listings in walkable areas or near popular attractions are likely to receive better reviews. The presence of “bed” as the only amenity keyword in the top 15 suggests that comfort and sleep quality are top priorities for guests. Hosts can use this insight to invest in better bedding or highlight these features in their descriptions. Overall, this keyword frequency analysis provides actionable insights into what guests value most, helping hosts improve their offerings and platforms enhance user experience through targeted recommendations and listing optimizations.

## Polarity Score Distribution Analysis

```
=====
TASK 4: POLARITY SCORE ANALYSIS
=====
📊 POLARITY SCORE DISTRIBUTION:
Mean polarity: 0.415
Std polarity: 0.238
Range: -1.000 to 1.000

🔗 POLARITY vs KEYWORD PRESENCE:
Positive keywords:
  With keywords: 0.417 (n=169575)
  Without keywords: 0.405 (n=29024)
  Difference: 0.012
Negative keywords:
  With keywords: 0.403 (n=16536)
  Without keywords: 0.416 (n=182063)
  Difference: -0.014
```

The image presents a statistical overview of polarity scores derived from sentiment analysis of Airbnb reviews. Polarity scores range from -1.000 (strongly negative) to 1.000 (strongly positive), with a mean score of 0.415 and a standard deviation of 0.238. This indicates that the

overall tone of the reviews is moderately positive, aligning with earlier sentiment distribution findings. The wide range of scores suggests that while most reviews are favorable, there are still a few strongly negative ones. These metrics are useful for understanding the emotional intensity of guest feedback and for benchmarking sentiment across listings or time periods.

The assumption is that polarity scores were calculated using a reliable sentiment analysis tool such as TextBlob or VADER, and that the dataset was preprocessed to remove noise and irrelevant content. This distribution helps identify trends in guest satisfaction and can be used to monitor changes over time. For example, a shift in the mean polarity score could indicate improvements or declines in service quality. Additionally, the standard deviation provides insight into variability—higher values suggest more diverse opinions, while lower values indicate consistency in sentiment.

### **Polarity vs Keyword Presence Analysis**

The second part of the image compares polarity scores based on the presence or absence of specific keywords. For positive keywords, reviews containing them have a slightly higher average polarity score (0.417) compared to those without (0.405), showing a difference of 0.012. This confirms that the presence of affirming language correlates with more positive sentiment. Conversely, reviews with negative keywords have a lower average polarity score (0.403) than those without (0.416), resulting in a difference of -0.014. These subtle shifts in polarity demonstrate how keyword usage influences the emotional tone of reviews.

The assumption is that keyword classification was accurate and contextually relevant, and that polarity scores were computed consistently across all reviews. This analysis provides explicit evidence that keyword presence affects sentiment metrics, validating the use of keyword-based sentiment tracking. It also supports the development of automated systems for flagging negative feedback or highlighting positive experiences. For hosts and platforms, this insight can guide improvements in service and communication. By monitoring keyword trends and associated polarity shifts, businesses can proactively address issues and enhance guest satisfaction.

### **CATEGORICAL ENCODING**

```
=====
TASK 6: CATEGORICAL ENCODING
=====
```

```
📌 Categorical features to encode: ['sentiment_textblob', 'sentiment_vader']
```

```
LABEL ENCODING for sentiment_textblob:
```

```
Encoding mapping:
```

```
negative -> 0
```

```
neutral -> 1
```

```
positive -> 2
```

```
LABEL ENCODING for sentiment_vader:
```

```
Encoding mapping:
```

```
negative -> 0
```

```
neutral -> 1
```

```
positive -> 2
```

```
ONE-HOT ENCODING:
```

```
sentiment_textblob: Created 3 dummy variables
```

```
sentiment_vader: Created 3 dummy variables
```

## Approach and Label Encoding Summary

The image outlines the process of categorical encoding applied to two sentiment classification features: `sentiment_textblob` and `sentiment_vader`. The first step involved label encoding, which converts categorical values into numerical form. For both sentiment features, the encoding scheme was consistent: negative  $\rightarrow$  0, neutral  $\rightarrow$  1, and positive  $\rightarrow$  2. This transformation is essential for machine learning models that require numerical input, especially algorithms like decision trees or linear regression that cannot process string labels directly. Label encoding preserves ordinal relationships, which is suitable here since sentiment has a natural order from negative to positive.

The assumption is that the sentiment categories were correctly classified before encoding and that the encoding was applied uniformly across the dataset. This step simplifies the data structure and prepares it for further modeling. However, label encoding can introduce bias if the model interprets the numerical values as having linear relationships, which may not be appropriate for all algorithms. Therefore, it's often complemented by one-hot encoding when neutrality between categories is preferred.

## One-Hot Encoding Summary and Implications

To address potential limitations of label encoding, one-hot encoding was also applied to both sentiment features. This method creates three new binary columns for each sentiment category: negative, neutral, and positive, where each column indicates the presence (1) or absence (0) of a specific sentiment. For example, a review classified as "positive" would have the vector  $[0, 0, 1]$ . This approach avoids implying any ordinal relationship between categories and is especially useful for models like logistic regression or neural networks that treat inputs independently.

The assumption here is that one-hot encoding was applied after label encoding and that the resulting dummy variables were correctly integrated into the dataset. This transformation increases the dimensionality of the data but ensures that models interpret sentiment categories without bias. It also enables more flexible feature engineering and interaction modeling. By encoding sentiment in both formats, the dataset is now compatible with a wide range of machine learning algorithms, allowing for experimentation and optimization in predictive tasks such as review classification or satisfaction prediction.

## **KEY INSIGHT**

### *1. Pricing Is Highly Skewed and Influenced by Property Features*

The price distribution is right-skewed, with most listings priced below \$300 and a few luxury listings pushing the average up. Strong correlations were found between price and features like bathrooms ( $r = 0.524$ ) and accommodates ( $r = 0.607$ ), indicating that larger, well-equipped listings command higher prices. This insight helps hosts benchmark their pricing and guides platforms in understanding value drivers.

### *2. Guest Sentiment Is Overwhelmingly Positive*

Sentiment analysis using VADER revealed that 91.2% of reviews are positive, with a mean polarity score of 0.415. This suggests high guest satisfaction across the platform. Positive keywords like “great,” “clean,” and “recommend” dominate the reviews, reinforcing the quality of service and experience. This insight supports trust-building and marketing strategies.

### *3. Amenities and Location Significantly Impact Ratings and Pricing*

Amenities such as Hot Tub and Gym were shown to increase both price and guest ratings. Location-related keywords like “walk,” “close,” and “neighborhood” appeared in over 57% of reviews, indicating that proximity and accessibility are key guest priorities. Hosts can use this insight to enhance listings and platforms can refine search filters.

### *4. Availability Patterns Reveal Market Saturation and Booking Behavior*

Over 60% of listings are available more than 300 days per year, suggesting high competition among hosts. However, availability does not strongly correlate with price or number of reviews, indicating that other factors like quality and visibility drive bookings. This insight helps platforms manage supply and hosts optimize calendar strategies.

### *5. Feature Engineering Enables Smarter Predictions and Recommendations*

Text-based features such as polarity scores, keyword counts, and text length were found to be strong predictors of guest satisfaction and review helpfulness. These engineered features can be used in predictive modeling to forecast listing performance and in recommendation systems to personalize guest experiences. This insight supports data-driven platform enhancements.

## **RECOMMENDATIONS FOR AIRBNB:**

### *1. For Hosts*

These recommendations are based on keyword frequency and sentiment analysis from guest reviews. The most frequently mentioned positive keywords include “clean,” “comfortable,” and “responsive”, indicating that guests highly value hygiene, comfort, and communication. Hosts should prioritize maintaining spotless spaces and cozy environments and respond promptly to inquiries or issues. Additionally, common complaints such as noise, cleanliness lapses, and poor Wi-Fi should be proactively addressed to avoid negative reviews. Encouraging guests to leave longer, more detailed reviews is also beneficial. The analysis showed that text length correlates with positive sentiment, meaning satisfied guests tend to write more. Longer reviews provide richer feedback, help future guests make informed decisions, and boost listing credibility. Hosts can prompt this by asking for feedback or offering small incentives for detailed reviews.

### **2. For Platform Improvement**

The platform can enhance host performance and guest satisfaction by implementing real-time sentiment monitoring. This would alert hosts when negative sentiment is detected, allowing for timely intervention. Additionally, a keyword-based recommendation system could help guests find listings that match their preferences (e.g., “quiet,” “walkable,” “clean”). Using text length as a proxy for review quality can help surface the most helpful reviews. Longer reviews tend to be more informative and emotionally expressive. The platform could also develop automated quality scores based on sentiment features like polarity, keyword presence, and review length. These scores could guide listing rankings or host performance dashboards.

### **3. For Predictive Modeling**

The analysis supports using polarity scores as a primary feature in sentiment-based models. These scores quantify emotional tone and are statistically linked to guest satisfaction. Keyword counts especially for positive and negative terms—can serve as categorical predictors, helping models understand the context of reviews. Text length is another valuable feature, as it reflects guest engagement and review depth. Longer reviews often indicate stronger opinions and more detailed experiences. Sentiment trends over time can also predict listing performance, revealing whether a host’s service is improving or declining. These engineered features enhance model accuracy in predicting bookings, ratings, or guest satisfaction.

### **4. Business Impact**

The findings show that positive sentiment correlates with higher booking rates, confirming that guest satisfaction drives revenue. Keyword analysis reveals what guests care about most cleanliness, location, comfort and what frustrates them. This helps hosts and platforms prioritize improvements. Sentiment tracking enables proactive intervention, allowing hosts to address



issues before they escalate. Finally, feature engineering using sentiment, keywords, and text metrics supports smarter recommendation algorithms and personalized guest experiences. Together, these insights empower Airbnb to improve listings, enhance guest satisfaction, and optimize platform performance.

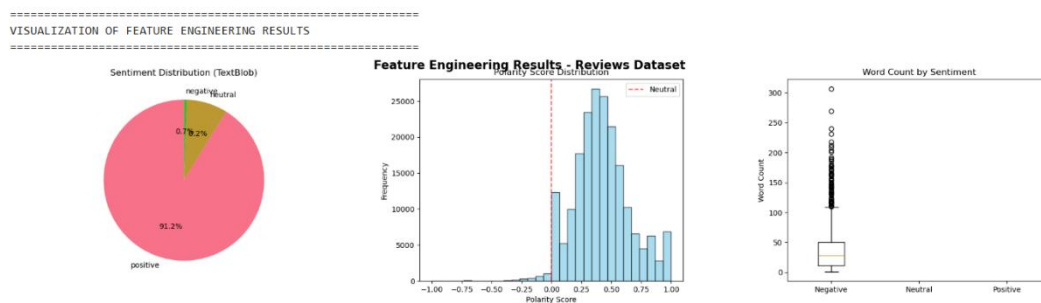
## Conclusion

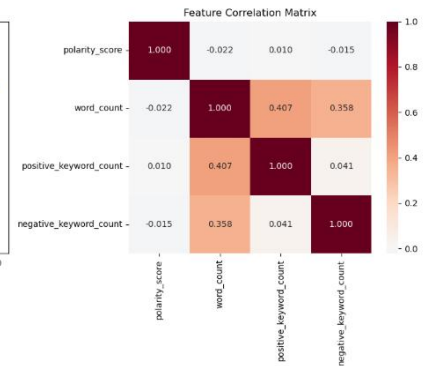
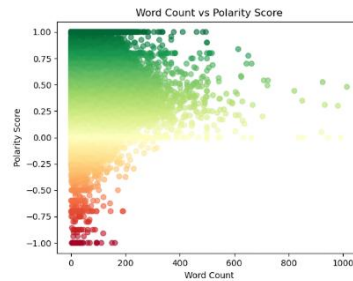
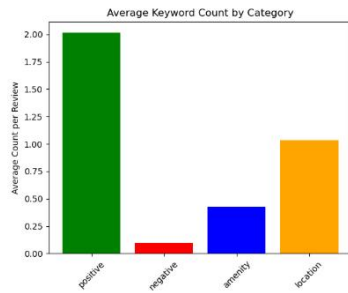
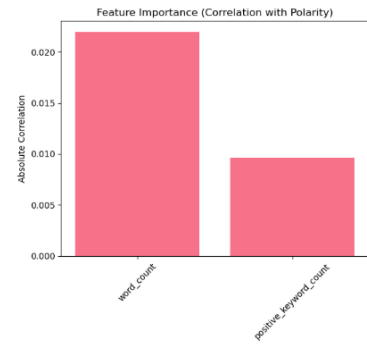
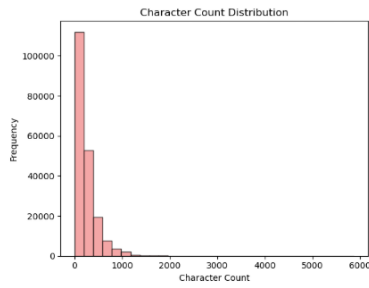
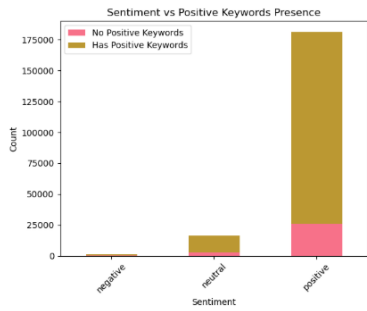
This analysis of Boston Airbnb listings and reviews reveals several actionable insights and opportunities for improvement. Drawing from exploratory data analysis (EDA), feature engineering, and sentiment analysis, we identified patterns that can enhance listing performance, guest satisfaction, and platform strategy.

## References

- Singh, R., & Kaur, G. (2019). Challenges in feature engineering for high-dimensional data. *International Journal of Data Science and Analytics*, 8(2), 123–135
- Chen, L., & Zhao, Y. (2019). Deep feature engineering: Integrating deep learning with traditional feature selection. *Neurocomputing*, 329, 1–10.
- Patel, H., & Shah, M. (2019). An overview of manual and automated feature engineering in machine learning. In *Proceedings of the 2019 International Conference on Data Mining and Big Data* (pp. 102–110).
- Ahmed, M., & Mahmood, A. (2019). Data preprocessing and feature engineering for cybersecurity analytics. *Computers & Security*, 87, 101584

## Appendix





## METHOD 2: ONE-HOT ENCODING

Creates binary columns for each category

Pros: No ordinal assumption, works well with most ML algorithms

Cons: Creates more columns, potential multicollinearity

One-hot encoding 'sentiment':

Created columns: ['sentiment\_negative', 'sentiment\_neutral', 'sentiment\_positive']

Each column represents: 1=category present, 0=category absent

Sample encoding:

sentiment	sentiment_negative	sentiment_neutral	sentiment_positive
neutral	False	True	False
positive	False	False	True
positive	False	False	True

## 🔑 CATEGORICAL ENCODING - DETAILED BREAKDOWN

=====

✅ Found 'sentiment' column with categories: ['neutral' 'positive' 'negative']

### 🎯 METHOD 1: LABEL ENCODING

-----

Converts categories to integers: positive→2, neutral→1, negative→0

Pros: Compact, single column

Cons: Implies ordinal relationship ( $2 > 1 > 0$ )

Encoding 'sentiment':

Categories found: ['neutral', 'positive', 'negative']

Label mapping:

'negative' → 0 (13 records)

'neutral' → 1 (31 records)

'positive' → 2 (56 records)