

MEMBERS & ROLES

Team Member	Role
Sai Akshith Bandari	Analytical Lead
Leslie Buckholtz	PM & Conceptual Design Lead
<u>Cibi Siddarth</u>	Technical Lead

1. Data Description with Tables

The dataset is a comprehensive collection of AirBnB property listings enriched with detailed metadata. It includes key identifiers such as listing and host IDs, which uniquely distinguish each property and its associated host. The dataset also captures various property features like capacity (accommodates, bedrooms, beds, bathrooms), descriptive textual data regarding the bathrooms, and categorical details such as property type and room type. Additionally, an extensive list of amenities offered per listing is provided, illuminating the various features that make each property unique. This detailed account ensures that every variable is well-defined, allowing for robust analysis of market trends, host performance, and guest preferences.

Tables	Parameters	Table Names	Description							
Timeline	last_scraped host_since calendar_last_scraped first_review last_review date_diff	tbl_Timeframe	Description	Listing ID	last_scraped	host_since	calendar_last_scraped	first_review	last_review	Date Diff First and Last Review
			Mean			5/30/2017		3/21/2021	4/28/2024	2.569025385
			Standard Error					5.645197165	2.016556574	0.013718933
			Median		9/21/2024	8/19/2016	9/21/2024	3/19/2022	8/7/2024	1.44109589
			Mode		12/23/2024	7/22/2012	12/23/2024	5/27/2024	9/2/2024	0
			Standard Deviation					1095.531927	391.3418865	2.927409447
			Sample Variance					1200190.202	153148.4721	8.569726073
			Kurtosis					-0.03293072	25.2453725	0.645814775
			Skewness					-0.89602129	-4.39634381	1.204631147
			Range					6096	5031	16.33972603
			Minimum		6/24/2024	3/3/2008	6/24/2024	6/22/2008	5/23/2011	0
			Maximum		3/2/2025	2/19/2025	3/2/2025	3/1/2025	3/1/2025	16.33972603
			Count	21,208	45,534		45,534	37,661	37,661	45,533
			Confidence Level(95.0%)					11.06473874	3.952505288	0.02688933

## Milestone Report 2: Data Processing and Exploratory Data Analysis

Occupancy Rates	availability_30	tbl_Occupancy	estimated_occupancy_l365d
	availability_60		65d
	availability_90		availability_30
	availability_365		availability_60
	availability_eoy		availability_90
	has_availability		availability_365
	estimated_occupancy_l365d		availability_eoy

## Milestone Report 2: Data Processing and Exploratory Data Analysis

Location	neighbourhood neighbourhood_cleansed neighbourhood_group_cleansed host_location host_neighbourhood latitude longitude	tbl_Location	<b>Most represented neighborhoods</b>						
			Mission Beach	4192					
			Pacific Beach	3562					
			La Jolla	2181					
			North Park	1525					
			<b>Least represented neighborhoods</b>						
			San Mateo Park	1					
			Fraccionamiento Torres Del Lago	1					
			Forestland	1					
Flamingos	1								
Five Points	1								
Property Type and Amenities	property_type room_type accommodates bathrooms bathrooms_text bedrooms beds amenities	tbl_PropertyType							
				accommodates	bathrooms	bedrooms	beds		
			Mean	4.747265779	1.432731585	1.824614574	2.340361049		
			Standard Error	0.015071361	0.004965295	0.006639071	0.010758373		
			Median	4	1	1	2		
			Mode	2	1	1	1		
			Standard Deviation	3.216032024	1.059529348	1.416691243	2.295696725		
			Sample Variance	10.34286198	1.122602439	2.007014077	5.270223453		
			Kurtosis	2.044409551	22.78499607	189.0568688	86.05815525		
			Skewness	1.415135755	2.688836282	5.455964523	4.672849646		
Range	15	23	66	85					
Minimum	1	0	0	0					
Maximum	16	23	66	85					
Confidence Level(95.0%)	0.029540109	0.009732059	0.013012686	0.021086585					
Seasonality	calendar_updated calendar_last_scraped availability_30, availability_60, availability_90, availability_365 first_review, last_review minimum_nights_avg_ntm, maximum_nights_avg_ntm	tbl_Seasonality							
				minimum_nights_avg_ntm	maximum_nights_avg_ntm	availability_30	availability_60	availability_90	availability_365
			Mean	13.21918717	913433.4781	12.97373391	30.86219089	51.20599991	199.9913691
			Standard Error	0.132072099	205696.474	0.050324896	0.098697982	0.143800704	0.57004211
			Median	3	365	12	32	57	212
			Mode	2	1125	0	0	0	0
			Standard Deviation	28.17813253	43886199.6	10.7386772	21.06086354	30.6851969	121.6395601
			Sample Variance	794.0071527	1.926E+15	115.319188	443.5599729	941.5813085	14796.18258
			Kurtosis	267.4460235	2373.803761	-1.371950137	-1.369655502	-1.169283847	-1.342166318
			Skewness	12.33539757	48.68180425	0.236568641	-0.127584086	-0.394206419	-0.189040288
Range	1124	2147483646	30	60	90	365			
Minimum	1	1	0	0	0	0			
Maximum	1125	2147483647	30	60	90	365			
Sum	601737.4	41579491923	590746	1405279	2331614	9106407			
Count	45520	45520	45534	45534	45534	45534			
Confidence Level(95.0%)	0.258863441	403168.4012	0.098637606	0.193449632	0.281851693	1.117291706			

## Milestone Report 2: Data Processing and Exploratory Data Analysis

Host Responsiveness	host_response_time host_response_rate host_acceptance_rate host_is_superuser	tbl_Host	host_response_rate						
			Mean	0.971278882					
			Standard Error	0.000595694					
			Median	1					
			Mode	1					
			Standard Deviation	0.120406857					
			Sample Variance	0.014497811					
			Kurtosis	44.14388068					
			Skewness	-6.355740333					
			Range	1					
			Minimum	0					
			Maximum	1					
			Count	40856					
			Confidence Level(95.0%)	0.001167574					
Market Demand	host_listings_count host_total_listings_count calculated_host_listings_count calculated_host_listings_count_entire_homes calculated_host_listings_count_private_rooms calculated_host_listings_count_shared_rooms estimated_occupancy_l365d estimated_revenue_l365d	tbl_MktDemand		estimated_occupancy_l365d	estimated_revenue_l365d	host_listings_count	host_total_listings_count	calculated_host_listings_count	calculated_host_listings_count_entire_homes
			Mean	14.05681469	17862.11693	112.5316389	165.501671	16.78793868	13.17870163
			Standard Error	0.238974981	331.7791263	2.575797747	3.554649682	0.16080069	0.129186651
			Median	0	8431.5	4	6	3	2
			Mode	0	0	1	1	1	1
			Standard Deviation	50.99414763	25725.18728	549.3277468	758.082696	34.31277251	27.56674831
			Sample Variance	2600.403093	661785260.6	301760.9734	574689.3739	1177.366357	759.9256123
			Kurtosis	14.16365909	18.96642682	51.756702	42.57526367	19.44017357	13.97050391
			Skewness	3.884479569	3.179565576	6.963972468	6.353016696	3.922410727	3.451462789
			Range	255	330225	5254	9054	263	178
			Minimum	0	0	1	1	1	0
			Maximum	255	330225	5255	9055	264	178
			Sum	640063	107387047	5118164	7527347	764422	600079
			Count	45534	6012	45482	45482	45534	45534
Confidence Level(95.0%)		0.468394806	650.4061026	5.048605172	6.967170769	0.31517194	0.253207914		

### 2. Key findings from your exploratory data analysis.

A comprehensive exploratory data analysis (EDA) was conducted to uncover meaningful insights, trends, anomalies, and relationships within the Airbnb dataset. The goal of this analysis is to guide the development of predictive and clustering models while informing host strategy and business decisions. Key findings are summarized below:

#### 2.1 Temporal and Group-Based Trends

- **Seasonal Demand Patterns:** Listings experience a clear seasonal trend, with review counts peaking between June and August, indicating increased guest activity during the summer. This suggests opportunities for dynamic pricing and targeted marketing during high-demand periods.
- **Listing Growth Over Time:** A year-over-year increase in total listings was observed, particularly post-2019. This growth aligns with Airbnb's rising popularity and may suggest an increasingly competitive market landscape in recent years.
- **Superhost Advantage:** Superhosts consistently outperform regular hosts across key metrics. They typically:
  - Charge 10–20% higher nightly rates
  - Receive more frequent and higher-rated reviews
  - Maintain lower cancellation rates and faster response times

These findings emphasize the value of achieving Super Host status to maximize occupancy and revenue.

#### 2.2 Outliers and Anomalies

- **Extreme Pricing:** A small subset of listings priced over \$1,000 per night significantly skew the overall price distribution. These listings likely represent luxury accommodations or potential data entry errors. Box plots and histograms helped identify and flag these outliers for further cleaning or segmentation.
- **Zero-Review Anomalies:** Numerous listings show full-year availability (365 days) but have zero guest reviews, which may indicate inactive or placeholder listings. These anomalies can distort average availability and occupancy calculations and should be filtered in modeling.
- **Review Count Extremes:** While the average listing has ~30–40 reviews, some listings exceed 500+ reviews, indicating highly established or corporate-managed properties. These high-volume listings require normalization when evaluating typical host performance.

### 2.3 Variable Relationships and Correlations

- Price vs. Review Count: Listings priced between \$100–\$200 generally receive the most reviews, suggesting this is the market’s “sweet spot.” Both ultra-low and high-priced listings tend to have lower engagement.
- Availability vs. Popularity: Listings with higher availability show a strong positive correlation with review count and guest feedback, implying that being available year-round increases occupancy potential.
- Review Score Breakdown:
- Cleanliness, Communication, and Location ratings have the highest correlation with overall score (Pearson’s  $r > 0.75$ ).
- These three factors should be considered critical KPIs (Key Performance Indicators) for hosts looking to improve listing quality and ratings.
- Room Type Impact: “Entire home/apartment” listings consistently outperform private or shared rooms in both price and rating, indicating a clear guest preference. These room types command higher average nightly prices and contribute to higher guest satisfaction.

### 2.4 Preliminary Sentiment Insights from Guest Reviews

- An initial sentiment analysis of guest reviews revealed overwhelmingly positive sentiment, with common keywords including:
- “Clean,” “Convenient location,” “Responsive host,” “Comfortable stay,” and “Easy check-in.”
- These terms reflect guest priorities and should be highlighted in listings to attract more bookings.
- Full sentiment scoring (e.g., using TextBlob or VADER) and word cloud visualization will be included in later stages to identify themes that contribute most to 5-star reviews.

## Conclusion

The EDA provided valuable insights into pricing dynamics, seasonal demand, host performance, and guest expectations. These findings lay the groundwork for selecting the most appropriate quantitative techniques in the next phase of the project, including:

- Regression models for price prediction
- Clustering algorithms for identifying high-performing neighborhoods
- Time series analysis for seasonality and occupancy forecasting
- Sentiment analysis to link guest feedback with listing performance

3. Propose and justify potential methods for your next analysis or modeling.

Based on the exploratory data analysis (EDA) and understanding of the Airbnb dataset, we propose a multi-faceted modeling approach tailored to address the core business objectives: accurate price prediction, neighborhood performance clustering, and actionable host insights through visualizations.

	Proposed Method	Justification
1.Accurate Predictions: Achieve high accuracy in predicting listing prices using machine learning models.	Multiple Linear Regression Random Forest Regressor XGBoost Regressor Artificial Neural Networks (ANN)	These models can effectively handle complex relationships in pricing data. Treebased models (Random Forest, XGBoost) are robust to outliers and non-linearities. ANN captures complex feature interactions.
2.Effective Clustering: Successfully identify high-performing neighborhoods through clustering analysis.	K-Means Clustering DBSCAN Hierarchical Clustering	Clustering helps group neighborhoods based on metrics like price, review scores, and availability. K-Means is effective for structured groups; DBSCAN detects noise/outliers; Hierarchical shows nested clusters.
3.Insightful		
Visualizations: Create clear and informative visualizations of geographic trends.	POWER BI DASHBOARD	
	<b>Visualizations for General Information Purposes</b>  <b>1. Price Distribution:</b> <u>A histogram or box plot</u> showing the distribution of listing prices. This can help users understand the range and average prices in different neighborhoods <b>2. Availability Heatmap:</b> <u>A heatmap</u> displaying the availability of listings throughout the year. This can help users identify peak and off-peak seasons <b>3. Geographic Trends:</b> <u>A map visualization</u> showing the geographic distribution of listings and their prices. This can help users visualize trends and identify high-demand areas <b>4. Review Sentiment Analysis:</b> A word cloud or sentiment score chart showing the most common words and overall sentiment in guest reviews. This can help users	<b>Visualizations to Help a User Become a Super Host</b>  <b>1. Host Response Time:</b> <u>A bar chart</u> showing the average response time of hosts. Faster response times can lead to better guest experiences and higher ratings <b>2. Review Scores Breakdown:</b> <u>A detailed breakdown of review scores</u> for accuracy, cleanliness, check-in, communication, location, and value. This can help hosts identify specific areas to improve <b>3. Monthly Review Trends:</b> <u>A line chart</u> showing the number of reviews per month. This can help hosts track their performance over time and identify trends <b>4. Amenities Analysis:</b> <u>A bar chart</u> showing the frequency of different amenities offered by super hosts. This can help hosts understand which amenities are most valued by guests



4. Sentiment Analysis: Provide meaningful insights from guest reviews and their impact on listing success.	understand guest feedback and areas for improvement	
	5. Neighborhood Clusters: A scatter plot or cluster map showing different neighborhoods based on rental performance and availability. This can help users identify high-performing neighborhoods	5. Occupancy and Revenue Trends: A line chart showing the estimated occupancy and revenue over the past year. This can help hosts understand their financial performance and identify opportunities for growth
	<ul style="list-style-type: none"><li>TextBlob or VADER for sentiment scoring</li><li>TF-IDF or CountVectorizer for keyword extraction</li><li>WordCloud for visual sentiment representation</li></ul>	Using Natural Language Processing (NLP), we will extract sentiment scores from guest reviews and analyze their impact on: <ul style="list-style-type: none"><li>Price variation</li><li>Overall Ratings</li><li>Booking frequency</li></ul>

Conclusion:

The proposed methods align with the key objectives of the project:

- Predictive models will optimize pricing strategy.
- Clustering will uncover high-performing neighborhoods.
- Power BI visualizations will provide actionable insights for hosts and business users.
- Sentiment analysis will connect guest feedback with listing performance and success metrics.