

## Exploratory Data Analysis (EDA) for Capstone

Name: Hosain Ahmed

Student ID: 301209637

Data Source: <https://www.kaggle.com/dgomonov/new-york-city-airbnb-open-data>

Data Structure:

- Rows: 48895
- Columns: 16
- Column names: id, name, host\_id, host\_name, neighbourhood\_group, neighbourhood, latitude, longitude, room\_type, price, minimum\_nights, number\_of\_reviews, last\_review, reviews\_per\_month, calculated\_host\_listings\_count, availability\_365
- Data types:
  - Numerical: 10
  - Categorical: 6
- Data Quality : EDA Report Link: [file:///D:/Capstone/eda\\_report\\_hosain\\_ahmed.html](file:///D:/Capstone/eda_report_hosain_ahmed.html)



## Overview

Dataset Statistics	
Number of Variables	16
Number of Rows	48895
Missing Cells	20141
Missing Cells (%)	2.6%
Duplicate Rows	0
Duplicate Rows (%)	0.0%
Total Size in Memory	23.5 MB
Average Row Size in Memory	504.1 B
Variable Types	Numerical: 10 Categorical: 6

## Exploratory Data Analysis (EDA) for Capstone

Name: Hosain Ahmed

Student ID: 301209637

### Dataset Insights

<code>last_review</code> has 10052 (20.56%) missing values	Missing
<code>reviews_per_month</code> has 10052 (20.56%) missing values	Missing
<code>host_id</code> is skewed	Skewed
<code>longitude</code> is skewed	Skewed
<code>price</code> is skewed	Skewed
<code>minimum_nights</code> is skewed	Skewed
<code>number_of_reviews</code> is skewed	Skewed
<code>reviews_per_month</code> is skewed	Skewed
<code>calculated_host_listings_count</code> is skewed	Skewed
<code>availability_365</code> is skewed	Skewed

1 2

### Dataset Insights

<code>name</code> has a high cardinality: 47905 distinct values	High Cardinality
<code>host_name</code> has a high cardinality: 11452 distinct values	High Cardinality
<code>neighbourhood</code> has a high cardinality: 221 distinct values	High Cardinality
<code>last_review</code> has a high cardinality: 1764 distinct values	High Cardinality
<code>last_review</code> has constant length 10	Constant Length
<code>longitude</code> has 48895 (100.0%) negatives	Negatives
<code>number_of_reviews</code> has 10052 (20.56%) zeros	Zeros
<code>availability_365</code> has 17533 (35.86%) zeros	Zeros

1 2

## Exploratory Data Analysis (EDA) for Capstone

Name: Hosain Ahmed

Student ID: 301209637

### Excluded Variables:

EXCLUSION	VARIABLE	DEFINITION
Irrelevant variables	id	Property ID
	Name	Property Tag Name
	host_id	Property host ID
	last_review	When the last time host recieved a review
Variables with more skewness value of 7.56774264	calculated_host_listings_count	Number of property a single host is offering

### Selected Variables:

VARIABLE	DEFINITION
<b>neighbourhood_group</b>	Larger neighborhood groups
<b>neighbourhood</b>	Specific áreas under the neighborhood groups
<b>latitude</b>	Latitude of the property
<b>longitude</b>	Longitude of the property
<b>room_type</b>	Types of rooms offered
<b>price</b>	Price per night stay in the property.
<b>minimum_nights</b>	Number of night stayed in the property
<b>number_of_reviews</b>	Total number of reviews each host received
<b>reviews_per_month</b>	Reviews received per month by each host
<b>availability_365</b>	The property was available for booking for how many days during the whole year

### Issues Found in variables:

Issues Found	Variables	Potential Actions
Skewness	host_id (Rejected), longitude, price, minimum_nights,	Cap & Floor and Log Transformation.

## Exploratory Data Analysis (EDA) for Capstone

Name: Hosain Ahmed

Student ID: 301209637

	number_of_reviews, reviews_per_month, calculated_host_listings_count, availability_365	
Missing Values	last_review, reviews_per_month	Imputation/ Drop Missing Values

- Potential actions
  - Cap & Floor for Skewness
  - Imputation for missing values
  - Log Transform for Skewness