TT  **B**  *I*  <>  🔗  🖼  99  ⅔≣  ☰  —  Ψ  ☺  ┄

```
#Project- <font color="ORANGE">Part A</font>: Airbnb Price Predictio
Insights:
Your Name: Bhushan Sahu
Course: [Your Course/Project Name]
batch : 15 january 2025
Project Goal : The goal of this project was to build a model to pred
listing prices using the 'Airbnb_Open_Data' dataset."
Dataset: "This dataset contains information on various Airbnb listin
including price, location, amenities, and host details."
```

# Project- <span style="color:orange">Part A</span>: Airbnb Price Prediction and Insights:

Your Name: Bhushan Sahu Course: [Your Course/Project Name] batch : 15 january 2025 Project Goal : The goal of this project was to build a model to predict Airbnb listing prices using the 'Airbnb_Open_Data' dataset." Dataset: "This dataset contains information on various Airbnb listings, including price, location, amenities, and host details."

## ❯ Load data

### Subtask:

Load the dataset from the CSV file into a pandas DataFrame.

[ ] ↳ 2 cells hidden

## ⌄ Explore data

### Subtask:

Analyze the dataset for trends, missing values, and outliers.

**Reasoning**: Print the shape and display the first few rows of the DataFrame as requested in the instructions.

```
print("Shape of the DataFrame:", df.shape)
display(df.head())
```

Shape of the DataFrame: (13251, 29)

| | id | log_price | property_type | room_type | amenities | accommodates | bathrooms | bed_type | cancellation_policy | cleanir |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6901257 | 5.010635 | Apartment | Entire home/apt | {"Wireless Internet","Air conditioning",Kitche... | 3.0 | 1.0 | Real Bed | strict | |
| 1 | 6304928 | 5.129899 | Apartment | Entire home/apt | {"Wireless Internet","Air conditioning",Kitche... | 7.0 | 1.0 | Real Bed | strict | |
| 2 | 7919400 | 4.976734 | Apartment | Entire home/apt | {TV,"Cable TV","Wireless Internet","Air condit... | 5.0 | 1.0 | Real Bed | moderate | |
| 3 | 13418779 | 6.620073 | House | Entire home/apt | {TV,"Cable TV",Internet,"Wireless Internet",Ki... | 4.0 | 1.0 | Real Bed | flexible | |
| 4 | 3808709 | 4.744932 | Apartment | Entire home/apt | {TV,Internet,"Wireless Internet","Air conditio... | 2.0 | 1.0 | Real Bed | moderate | |

5 rows × 29 columns

**Reasoning**: Display the data types of each column and generate descriptive statistics for numerical columns to understand the data distribution and identify potential issues.

```
print("\nData types of each column:")
display(df.info())

print("\nDescriptive statistics for numerical columns:")
display(df.describe())
```

```
Data types of each column:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13251 entries, 0 to 13250
Data columns (total 29 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   id                     13251 non-null  int64
 1   log_price              13251 non-null  float64
 2   property_type          13251 non-null  object
 3   room_type              13251 non-null  object
 4   amenities              13251 non-null  object
 5   accommodates           13250 non-null  float64
 6   bathrooms              13210 non-null  float64
 7   bed_type               13250 non-null  object
 8   cancellation_policy    13250 non-null  object
 9   cleaning_fee           13250 non-null  object
 10  city                   13250 non-null  object
 11  description            13250 non-null  object
 12  first_review           10418 non-null  object
 13  host_has_profile_pic   13214 non-null  object
 14  host_identity_verified 13214 non-null  object
 15  host_response_rate     9965 non-null   object
 16  host_since             13214 non-null  object
 17  instant_bookable       13250 non-null  object
 18  last_review            10425 non-null  object
 19  latitude               13250 non-null  float64
 20  longitude              13250 non-null  float64
 21  name                   13250 non-null  object
 22  neighbourhood          12044 non-null  object
 23  number_of_reviews      13250 non-null  float64
 24  review_scores_rating   10261 non-null  float64
 25  thumbnail_url          11732 non-null  object
 26  zipcode                13077 non-null  object
 27  bedrooms               13238 non-null  float64
 28  beds                   13229 non-null  float64
dtypes: float64(9), int64(1), object(19)
memory usage: 2.9+ MB
None
```

Descriptive statistics for numerical columns:

|       | id | log_price | accommodates | bathrooms | latitude | longitude | number_of_reviews | review_scores_rating |
|-------|------|-----------|--------------|-----------|----------|-----------|-------------------|----------------------|
| count | 1.325100e+04 | 13251.000000 | 13250.000000 | 13210.000000 | 13250.000000 | 13250.000000 | 13250.000000 | 10261.000000 |
| mean  | 1.126456e+07 | 4.782862 | 3.142113 | 1.229334 | 38.441380 | -92.568106 | 21.021434 | 93.985187 |
| std   | 6.093090e+06 | 0.718365 | 2.137676 | 0.568895 | 3.075839 | 21.768362 | 38.339510 | 8.148463 |
| min   | 3.362000e+03 | 0.000000 | 1.000000 | 0.000000 | 33.339327 | -122.508663 | 0.000000 | 20.000000 |
| 25%   | 6.243545e+06 | 4.317488 | 2.000000 | 1.000000 | 34.129031 | -118.344198 | 1.000000 | 92.000000 |
| 50%   | 1.230964e+07 | 4.709530 | 2.000000 | 1.000000 | 40.663021 | -76.998867 | 6.000000 | 96.000000 |
| 75%   | 1.641626e+07 | 5.217646 | 4.000000 | 1.000000 | 40.745656 | -73.954436 | 24.000000 | 100.000000 |
| max   | 2.122746e+07 | 7.598399 | 16.000000 | 8.000000 | 42.389907 | -70.991861 | 492.000000 | 100.000000 |

**Reasoning**: Check for missing values in each column and display the count or percentage of missing values to understand the extent of missing data.

```
print("\nMissing values in each column:")
display(df.isnull().sum())
print("\nPercentage of missing values in each column:")
display(df.isnull().sum() / len(df) * 100)
```

Missing values in each column:

|  | 0 |
|---|---|
| id | 0 |
| log_price | 0 |
| property_type | 0 |
| room_type | 0 |
| amenities | 0 |
| accommodates | 1 |
| bathrooms | 41 |
| bed_type | 1 |
| cancellation_policy | 1 |
| cleaning_fee | 1 |
| city | 1 |
| description | 1 |
| first_review | 2833 |
| host_has_profile_pic | 37 |
| host_identity_verified | 37 |
| host_response_rate | 3286 |
| host_since | 37 |
| instant_bookable | 1 |
| last_review | 2826 |
| latitude | 1 |
| longitude | 1 |
| name | 1 |
| neighbourhood | 1207 |
| number_of_reviews | 1 |
| review_scores_rating | 2990 |
| thumbnail_url | 1519 |
| zipcode | 174 |
| bedrooms | 13 |
| beds | 22 |

**dtype:** int64

Percentage of missing values in each column:

|  | 0 |
|---|---|
| id | 0.000000 |
| log_price | 0.000000 |
| property_type | 0.000000 |
| room_type | 0.000000 |
| amenities | 0.000000 |
| accommodates | 0.007547 |
| bathrooms | 0.309411 |
| bed_type | 0.007547 |
| cancellation_policy | 0.007547 |
| cleaning_fee | 0.007547 |
| city | 0.007547 |
| description | 0.007547 |
| first_review | 21.379519 |
| host_has_profile_pic | 0.279224 |
| host_identity_verified | 0.279224 |
| host_response_rate | 24.798128 |

|  |  |
|---|---|
| host_response_rate | 24.798128 |
| host_since | 0.279224 |
| instant_bookable | 0.007547 |
| last_review | 21.326692 |
| latitude | 0.007547 |
| longitude | 0.007547 |
| name | 0.007547 |
| neighbourhood | 9.108747 |
| number_of_reviews | 0.007547 |
| review_scores_rating | 22.564335 |
| thumbnail_url | 11.463286 |
| zipcode | 1.313108 |
| bedrooms | 0.098106 |
| beds | 0.166025 |

**dtype:** float64

◄                                                                                                                    ▶

**Reasoning**: Identify categorical columns and examine the unique values and their counts to detect inconsistencies or rare categories.

```
categorical_cols = df.select_dtypes(include='object').columns
print("\nUnique values and counts for categorical columns:")
for col in categorical_cols:
    print(f"\nColumn: {col}")
    display(df[col].value_counts())
```

|  |  |
|---|---|
| **Timeshare** | 12 |
| **Hostel** | 11 |
| **Boutique hotel** | 11 |
| **Boat** | 10 |
| **Cabin** | 10 |
| **Serviced apartment** | 9 |
| **Vacation home** | 2 |
| **Castle** | 1 |
| **Treehouse** | 1 |
| **Tipi** | 1 |
| **Tent** | 1 |
| **Hut** | 1 |

**dtype:** int64

Column: room_type

|  | count |
|---|---|
| **room_type** |  |
| **Entire home/apt** | 7402 |
| **Private room** | 5453 |
| **Shared room** | 396 |

**dtype:** int64

Column: amenities

|  | count |
|---|---|
| **amenities** |  |
| {} | 114 |
| {"translation missing: en.hosting_amenity_49","translation missing: en.hosting_amenity_50"} | 29 |
| {"Family/kid friendly"} | 18 |
| {TV,"Cable TV",Internet,"Wireless Internet","Air conditioning",Kitchen,"Pets allowed",Doorman,Gym,Elevator,Heating,"Family/kid friendly",Washer,Dryer,"Smoke detector","Carbon monoxide detector",Essentials,Shampoo,"24-hour check-in",Hangers,"Hair dryer",Iron,"Laptop friendly workspace","Self Check-In",Doorman} | 6 |
| {TV,"Cable TV",Internet,"Wireless Internet","Air conditioning",Kitchen,"Free parking on premises",Heating,"Family/kid friendly",Washer,Dryer,"Smoke detector","Carbon monoxide detector","First aid kit","Safety card","Fire extinguisher",Essentials,Shampoo,"24-hour check-in",Hangers,"Hair dryer",Iron,"Laptop friendly workspace"} | 6 |
| ... | ... |

| | |
|---|---|
| {TV,Internet,"Wireless Internet","Air conditioning",Kitchen,"Free parking on premises","Smoking allowed","Pets allowed",Heating,"Family/kid friendly","Smoke detector","First aid kit","Fire extinguisher","Lock on bedroom door",Hangers} | 1 |
| {Internet,"Wireless Internet","Air conditioning",Heating,"Smoke detector","Fire extinguisher",Essentials,Shampoo,Hangers,"Hair dryer",Iron,"translation missing: en.hosting_amenity_49","translation missing: en.hosting_amenity_50"} | 1 |
| {"Cable TV","Wireless Internet","Air conditioning","Wheelchair accessible",Kitchen,"Pets live on this property",Dog(s),"Buzzer/wireless intercom","Family/kid friendly","Smoke detector","Carbon monoxide detector","Fire extinguisher"} | 1 |
| {TV,"Cable TV",Internet,"Wireless Internet","Air conditioning",Kitchen,"Free parking on premises",Heating,Washer,Dryer,"Smoke detector","Safety card","Fire extinguisher",Essentials,Shampoo,"24-hour check-in",Hangers,"Hair dryer",Iron,"Laptop friendly workspace","translation missing: en.hosting_amenity_49","translation missing: en.hosting_amenity_50"} | 1 |
| {TV,"Cable TV",Internet,"Wireless Internet","Air conditioning",Kitchen,"Pets allowed",Heating,"Family/kid friendly","Suitable for events",Washer,Dryer,"Smoke detector","Carbon monoxide detector","First aid kit","Fire extinguisher",Essentials,Shampoo,"24-hour check-in",Hangers,"Hair dryer",Iron,"Laptop friendly workspace","translation missing: en.hosting_amenity_50"} | 1 |

12771 rows × 1 columns

**dtype:** int64

Column: bed_type

| | count |
|---|---|
| **bed_type** | |
| Real Bed | 12869 |
| Futon | 138 |
| Pull-out Sofa | 111 |
| Airbed | 78 |
| Couch | 54 |

**dtype:** int64

Column: cancellation_policy

| | count |
|---|---|
| **cancellation_policy** | |
| strict | 5822 |
| flexible | 3992 |
| moderate | 3416 |
| super_strict_30 | 16 |
| super_strict_60 | 4 |

**dtype:** int64

Column: cleaning_fee

| | count |
|---|---|
| **cleaning_fee** | |
| True | 9789 |
| False | 3461 |

**dtype:** int64

Column: city

| | count |
|---|---|
| **city** | |
| NYC | 5793 |
| LA | 4004 |
| SF | 1211 |
| DC | 976 |
| Chicago | 656 |
| Boston | 610 |

**dtype:** int64

Column: description

| | count |
|---|---|
| **description** | |
| A very cozy house located in the heart of Hollywood between two famous streets: Beverly Blvd and Melrose Ave, within a safe and quiet | |

| | |
|---|---|
| neighborhood, 30 minutes from LAX, 15 minutes from Downtown LA, 15 minutes from Griffith Observatory, 15 minutes from Hollywood Walk of Fame, 15 minutes from Universal Studios, 5 minutes from Paramount Pictures. There are lots of restaurants and shops nearby. My favorite is Osteria La Buca:) A great house in the central area of Los Angeles suitable both for travelers and business trips. It is fully equipped with all you need and freshly renovated. FREE PARKING:) | 3 |
| Newly renovated studio apartment in prime Miracle Mile location. Queen size bed. Couch. TV. Full cable and internet. Kitchenette with the essentials- dishes, cups, cookware, fridge, microwave, toaster oven, coffee maker and more! Linens, towels and bath accessories all provided. The apartment is right off of Wilshire Blvd. so getting around will be a breeze. | 2 |
| Our apartment is beautifully designed with luxurious features like hardwood flooring, nine foot ceilings and bay or floor to ceiling windows. Residents can walk to the Metro, shops & restaurants. This beautiful luxury 14 floor Hi-rise building offers many fine amenities including a sun deck, rooftop swimming pool, business center, fitness facility, and so much more. Within our apartment you can exceptional features such as hardwood flooring, nine foot ceilings, and bay or floor to ceiling windows. We are pleased to ensure the comfort of our guests by providing linens, washer and dryer, dishwasher, and high speed internet access. The Master Bedroom features a deluxe Queen bed with our indulgent custom linens, fluffy duvet, and plush pillows, two night tables with lamps, clock radio, dresser and very spacious walk-in closets. The secondary bedroom offers a queen sized bed as well as spacious closet. The Living Room includes a sofa with pullout be | 2 |
| Good for couples, adventurers, business travelers, and families. This apartment has 3 private lockable BR, each with 1 queen bed. Each of the 3 rooms has a separate listing and can be reserved separately for larger groups or by other guests. Common kitchen, bath, dining rm, living rm. Plenty of free on-street parking available. Urban location 7mi south of downtown, easy trip. Uber costs about $15 or about half that if ride-share can be used. uber to the train station is about $5 then $2.25 ea each Bedroom has a keyless lock that can be locked from the outside for extra security. Guests have access to a private lockable bedroom. In addition, guests in the 3 bedrooms share a common bathroom, kitchen, living room and dining room. This is a second floor unit. Free street parking available | 2 |
| Upgraded 2bd 2bath Apartment in LA's most desirable location. Resort Style Swimming Pool. 2 Free Parking Spaces. Washer and Dryer in the unit. Free High Speed Internet. Luxurious modern apartment in the heart of LA The building is VERY secured with 24/7 security guards on premises 2 covered parking spaces for your cars plus guest parking 2 master bedrooms separated by the living room The apartment is perfect for two couples, a family, or two co-workers Each master bedroom has its own walk-in closet and large bathroom WASHER & DRYER in unit Fully equipped kitchen The balcony seats 4 people comfortably and has an electric BBQ 50" LED LCD TV with Apple TV and all the news channels High speed Internet Access to the large swimming pool area where you can swim and relax I will meet you at the apartment to check-you-in and give you the keys and go over all the details. Walk through first-class farmers' market and the Grove, contemplate art at the LACMA, and marvel at your central lo | 2 |
| ... | ... |
| Personal Use of Large backyard for event hosting and parties between 2-49 people. Includes use of treehouse by 4 people or less at a time, along with a Flat Roof for relaxing or star gazing, that's connected to the treehouse. This listing is for events or parties between 2- 49 persons interested in throwing events in a Backyard space. The listing includes a Treehouse built on an orange tree and a . chill spot underneath treehouse with couch and desk. Backyard also has a small table and 4 chairs for guest to use. laundry is in backyard shed. For parties over 5 people, or parties lasting over 2 days, a porta potty will need to be rented. This can be paid for directly or as agreed upon by me under the reservation for an extra fee of $110.00. Treehouse, with chill spot underneath. Roof from treehouse. Entire backyard space. Coin laundry facilities outside if necessary. Hose and running water outside for water play or portable pool and/or water slide play. I will share anything I have to h | 1 |
| My place is a Penthouse with 3 Bedrooms and 3 Bathrooms with wrap around views of the city and access to a private big rooftop, and is centrally located on a quiet street, in a doorman building and just two blocks from all major subway lines!. Is close to South Street Seaport, Wall St, and New York Stock Exchange. You'll love my place, the outdoor space, the light, the neighborhood, and the ambiance. My place is good for couples, solo adventurers, and business travelers. One 65' 4k TV, available computers. 300 mb internet, Netflix and prime streaming. If am available, I will interact with the guests as much as they will like me too. | 1 |
| Spacious DC Condo built in 2004 with a 400 sq ft balcony, located in the heart of downtown DC. 1 block from Metro, Verizon Center. 6 blocks to National Mall, walking distance to Capitol, White House, Smithsonian. 1BR + a den that is perfect for a family with young children. | 1 |
| Spacious modern loft living in the heart of Downtown Brooklyn. Real artist loft in the middle of the action. The is the real deal. Large, open and airy loft on the 5th floor in a former factory building. It has one bedroom with a queen-sized bed, a large kitchen, bath with shower, television, wireless and views over Fort Greene and the Brooklyn Navy Yard. The building has part-time elevator service (11am-11pm M-F, 12-10pm S-S), so you may have to walk up in the late evenings and mornings. There are seven train lines in the area and most go to lower Manhattan in one or two stops. The trains are A/C/F/N/R/Q/B/G/4/5/2/3. The building is walking distance to the Brooklyn and Manhattan Bridges, DUMBO, Carroll Gardens, Fort Greene and the Nets Stadium. Cabs are easily found on Flatbush Avenue around the corner. The apartment is 1000 square feet and has 12 foot ceilings. There is also a shared roof deck with lovely Manhattan views. The couch can accommodate a third guest, if needed. This is do | 1 |
| A sunny room in a historic house in the heart of Park Slope, one block from Prospect Park, near Brooklyn Academy of Music, Brooklyn Museum, the Brooklyn Botanical Gardens, and Barclays Center. A block from stores and restaurants. Beautifully restored with oak floors, marble fireplace, brass lights, artwork on the walls, indoor plants, overlooking a garden. Central air, radiant heat. The bathroom has a skylight, walk-in shower, marble and cherry vanity, and Jerusalem Gold tiles. The house was built in 1885 and fully renovated in 2009 to a very high level. It is considered one of the nicest buildings in Park Slope. The guest bedroom and bathroom occupy part of the top floor and are linked via a hallway above which is a skylight. Stairs lead up from the ground floor in the center of the house. I've lived in Park Slope for 30+ years and know New York City well. Happy to advise on sights, experiences, and restaurants. As former competitive runner, I know good routes in the park ne | 1 |

13222 rows × 1 columns

**dtype:** int64

Column: first_review

| | count |
|---|---|
| first_review | |
| 01-01-2017 | 57 |
| 22-01-2017 | 43 |
| 03-01-2016 | 40 |
| 02-01-2017 | 40 |
| 04-09-2017 | 38 |
| ... | ... |

| | |
|---|---|
| **05-02-2015** | 1 |
| **24-09-2014** | 1 |
| **26-04-2013** | 1 |
| **19-08-2012** | 1 |
| **01-04-2012** | 1 |

1859 rows × 1 columns

**dtype:** int64

Column: host_has_profile_pic

| | count |
|---|---|
| **host_has_profile_pic** | |
| t | 13179 |
| f | 35 |

**dtype:** int64

Column: host_identity_verified

| | count |
|---|---|
| **host_identity_verified** | |
| t | 8968 |
| f | 4246 |

**dtype:** int64

Column: host_response_rate

| | count |
|---|---|
| **host_response_rate** | |
| **100%** | 7667 |
| **90%** | 437 |
| **80%** | 214 |
| **0%** | 162 |
| **50%** | 104 |
| **...** | ... |
| **35%** | 1 |
| **72%** | 1 |
| **6%** | 1 |
| **36%** | 1 |
| **13%** | 1 |

66 rows × 1 columns

**dtype:** int64

Column: host_since

| | count |
|---|---|
| **host_since** | |
| **30-03-2015** | 51 |
| **14-02-2014** | 31 |
| **29-07-2014** | 21 |
| **02-07-2014** | 19 |
| **16-05-2016** | 19 |
| **...** | ... |
| **30-12-2010** | 1 |
| **28-10-2010** | 1 |
| **13-01-2011** | 1 |
| **02-11-2009** | 1 |
| **05-06-2011** | 1 |

2649 rows × 1 columns

**dtype:** int64

Column: instant_bookable

|  | count |
| --- | --- |
| instant_bookable | |
| f | 9793 |
| t | 3457 |

**dtype:** int64

Column: last_review

|  | count |
| --- | --- |
| last_review | |
| 17-09-2017 | 238 |
| 30-04-2017 | 232 |
| 24-09-2017 | 215 |
| 23-04-2017 | 174 |
| 18-09-2017 | 149 |
| ... | ... |
| 01-04-2014 | 1 |
| 09-09-2016 | 1 |
| 11-02-2016 | 1 |
| 21-01-2014 | 1 |
| 06-02-2015 | 1 |

956 rows × 1 columns

**dtype:** int64

Column: name

|  | count |
| --- | --- |
| name | |
| Central Park Bliss | 2 |
| home away from home | 2 |
| Private Room Beautiful Mansion Beverly Hills | 2 |
| Modern Studio Apartment | 2 |
| Cozy studio | 2 |
| ... | ... |
| Manhattan on a Budget | 1 |
| Cozy, 1 bedroom Brownstone Apt | 1 |
| 1 min Walk to beach - Zen Beach Pad | 1 |
| Luxury 2BR Midtown East-Near UN! | 1 |
| Great bedroom in Mid City! | 1 |

13218 rows × 1 columns

**dtype:** int64

Column: neighbourhood

|  | count |
| --- | --- |
| neighbourhood | |
| Williamsburg | 518 |
| Bedford-Stuyvesant | 370 |
| Bushwick | 304 |
| Mid-Wilshire | 256 |
| Upper West Side | 250 |
| ... | ... |
| Bethesda, MD | 1 |

| | |
|---|---|
| **Grymes Hill** | 1 |
| **Albany Park** | 1 |
| **Arleta** | 1 |
| **San Marino** | 1 |

511 rows × 1 columns

**dtype:** int64

Column: thumbnail_url

| | count |
|---|---|
| **thumbnail_url** | |
| https://a0.muscache.com/im/pictures/fef753d9-a9a3-4109-81a0-ec434b3d867a.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/6d7cbbf7-c034-459c-bc82-6522c957627c.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/348a55fe-4b65-452a-b48a-bfecb3b58a66.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/6fae5362-9e3a-4fa9-aa54-bbd5ea26538d.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/226074de-610c-4fcd-a705-5828312cd261.jpg?aki_policy=small | 1 |
| ... | ... |
| https://a0.muscache.com/im/pictures/8d2f08ce-bf65-4018-a7b0-18823a7882a7.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/0ed6c128-7d60-4e05-b3bf-63158a230f70.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/61bd05d5-c4db-4c49-9f87-c0981c2d83b9.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/4c920c60-43dc-4169-a0da-ccf37f1d7a94.jpg?aki_policy=small | 1 |
| https://a0.muscache.com/im/pictures/82509143-4b21-44eb-a556-e3c1e0afac60.jpg?aki_policy=small | 1 |

11732 rows × 1 columns

**dtype:** int64

Column: zipcode

| | count |
|---|---|
| **zipcode** | |
| **11211** | 297 |
| **90291** | 234 |
| **11221** | 206 |
| **94110** | 198 |
| **10002** | 187 |
| **...** | ... |
| **11420** | 1 |
| **90604** | 1 |
| **94130** | 1 |
| **90061** | 1 |
| **60628** | 1 |

550 rows × 1 columns

**dtype:** int64

## Preprocess data

### Subtask:

Clean the data, perform feature engineering, and transformations.

**Reasoning**: Drop columns with a high percentage of missing values as identified in the previous exploration step.

```
columns_to_drop = ['host_response_rate', 'review_scores_rating', 'first_review', 'last_review', 'thumbnail_url', 'neighbourhood']
df = df.drop(columns=columns_to_drop)
display(df.head())
```

| | id | log_price | property_type | room_type | amenities | accommodates | bathrooms | bed_type | cancellation_policy | cleanin |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6901257 | 5.010635 | Apartment | Entire home/apt | {"Wireless Internet","Air conditioning",Kitche... | 3.0 | 1.0 | Real Bed | strict | |
| 1 | 6304928 | 5.129899 | Apartment | Entire home/apt | {"Wireless Internet","Air conditioning",Kitche... | 7.0 | 1.0 | Real Bed | strict | |
| 2 | 7919400 | 4.976734 | Apartment | Entire home/apt | {TV,"Cable TV","Wireless Internet","Air condit... | 5.0 | 1.0 | Real Bed | moderate | |
| 3 | 13418779 | 6.620073 | House | Entire home/apt | {TV,"Cable TV",Internet,"Wireless Internet",Ki... | 4.0 | 1.0 | Real Bed | flexible | |
| 4 | 3808709 | 4.744932 | Apartment | Entire home/apt | {TV,Internet,"Wireless Internet","Air conditio... | 2.0 | 1.0 | Real Bed | moderate | |

5 rows × 23 columns

**Reasoning**: Handle remaining missing values in numerical columns by filling with the mean or median. I will use the median as it is less sensitive to outliers.

```
numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
for col in numerical_cols:
    if df[col].isnull().sum() > 0:
        median_val = df[col].median()
        df[col].fillna(median_val, inplace=True)

display(df.isnull().sum())
```

| | 0 |
|---|---|
| id | 0 |
| log_price | 0 |
| property_type | 0 |
| room_type | 0 |
| amenities | 0 |
| accommodates | 0 |
| bathrooms | 0 |
| bed_type | 1 |
| cancellation_policy | 1 |
| cleaning_fee | 1 |
| city | 1 |
| description | 1 |
| host_has_profile_pic | 37 |
| host_identity_verified | 37 |
| host_since | 37 |
| instant_bookable | 1 |
| latitude | 0 |
| longitude | 0 |
| name | 1 |
| number_of_reviews | 0 |
| zipcode | 174 |
| bedrooms | 0 |
| beds | 0 |

**Reasoning**: Handle remaining missing values in categorical columns by filling with the mode or a placeholder like 'missing'. Also handle the remaining missing values in numerical columns that were missed in the previous step like 'zipcode'.

```python
categorical_cols = df.select_dtypes(include='object').columns
for col in categorical_cols:
    if df[col].isnull().sum() > 0:
        mode_val = df[col].mode()[0]
        df[col].fillna(mode_val, inplace=True)

# Fill remaining numerical missing values (like zipcode)
numerical_cols = df.select_dtypes(include=['float64', 'int64']).columns
for col in numerical_cols:
    if df[col].isnull().sum() > 0:
        median_val = df[col].median()
        df[col].fillna(median_val, inplace=True)

display(df.isnull().sum())
```

|  | 0 |
|---|---|
| id | 0 |
| log_price | 0 |
| property_type | 0 |
| room_type | 0 |
| amenities | 0 |
| accommodates | 0 |
| bathrooms | 0 |
| bed_type | 0 |
| cancellation_policy | 0 |
| cleaning_fee | 0 |
| city | 0 |
| description | 0 |
| host_has_profile_pic | 0 |
| host_identity_verified | 0 |
| host_since | 0 |
| instant_bookable | 0 |
| latitude | 0 |
| longitude | 0 |
| name | 0 |
| number_of_reviews | 0 |
| zipcode | 0 |
| bedrooms | 0 |
| beds | 0 |

**Reasoning**: Convert `host_since` to datetime objects and create a new feature 'host_age' representing the age of the host listing.

```
df['host_since'] = pd.to_datetime(df['host_since'], format='%d/%m/%Y', errors='coerce')
current_date = pd.to_datetime('today')
df['host_age'] = (current_date - df['host_since']).dt.days
display(df[['host_since', 'host_age']].head())
```

|  | host_since | host_age |
|---|---|---|
| 0 | NaT | NaN |
| 1 | NaT | NaN |
| 2 | NaT | NaN |
| 3 | NaT | NaN |
| 4 | NaT | NaN |

**Reasoning**: The previous attempt to convert `host_since` failed because the date format was incorrect. I will try converting it again with the correct format and then calculate the host age.

```
df['host_since'] = pd.to_datetime(df['host_since'], format='%d-%m-%Y', errors='coerce')
current_date = pd.to_datetime('2025-06-03')
df['host_age'] = (current_date - df['host_since']).dt.days
# Fill potential NaT values in host_since and resulting NaN in host_age after conversion
if df['host_since'].isnull().sum() > 0:
    # Fill with a placeholder date or the mode if appropriate, here using a placeholder
    df['host_since'] = df['host_since'].fillna(pd.to_datetime('2010-01-01'))
    # Recalculate host_age for the filled values
    df['host_age'] = (current_date - df['host_since']).dt.days

display(df[['host_since', 'host_age']].head())
display(df['host_age'].isnull().sum())
```

| | host_since | host_age |
|---|---|---|
| 0 | 2010-01-01 | 5632 |
| 1 | 2010-01-01 | 5632 |
| 2 | 2010-01-01 | 5632 |
| 3 | 2010-01-01 | 5632 |
| 4 | 2010-01-01 | 5632 |

**Reasoning**: Extract relevant information from the `amenities` column by creating dummy variables for the most common amenities.

```python
from collections import Counter
import json

# Clean and parse amenities
df['amenities'] = df['amenities'].str.replace('{', '').str.replace('}', '').str.replace('"', '')
all_amenities = Counter(amenity for amenities_list in df['amenities'] for amenity in amenities_list.split(','))

# Get the 50 most common amenities
most_common_amenities = [amenity for amenity, count in all_amenities.most_common(50)]

# Create dummy variables for the most common amenities
for amenity in most_common_amenities:
    # Handle cases where amenity is an empty string
    if amenity:
        df[amenity] = df['amenities'].apply(lambda x: 1 if amenity in x else 0)

display(df[most_common_amenities].head())
```

| | Wireless Internet | Kitchen | Heating | Essentials | Smoke detector | Air conditioning | TV | Shampoo | Hangers | Carbon monoxide detector | ... | Doorman | Wheelchair accessible | Refriger |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | ... | 0 | 0 | |
| 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | ... | 0 | 0 | |
| 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | ... | 0 | 0 | |
| 3 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | ... | 0 | 0 | |
| 4 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | ... | 0 | 0 | |

**Reasoning**: Create dummy variables for other categorical columns.

```python
categorical_cols_to_dummy = ['property_type', 'room_type', 'bed_type', 'cancellation_policy', 'cleaning_fee', 'host_has_profile_pic', 'H
df = pd.get_dummies(df, columns=categorical_cols_to_dummy, dummy_na=False)
display(df.head())
```

| | id | log_price | amenities | accommodates | bathrooms | description | host_since | latitude | longitude | name |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6901257 | 5.010635 | Wireless Internet,Air conditioning,Kitchen,Hea... | 3.0 | 1.0 | Beautiful, sunlit brownstone 1-bedroom in the ... | 2010-01-01 | 40.696524 | -73.991617 | Beautiful brownstone 1-bedroom |
| 1 | 6304928 | 5.129899 | Wireless Internet,Air conditioning,Kitchen,Hea... | 7.0 | 1.0 | Enjoy travelling during your stay in Manhattan... | 2010-01-01 | 40.766115 | -73.989040 | Superb 3BR Apt Located Near Times Square |
| 2 | 7919400 | 4.976734 | TV,Cable TV,Wireless Internet,Air conditioning... | 5.0 | 1.0 | The Oasis comes complete with a full backyard ... | 2010-01-01 | 40.808110 | -73.943756 | The Garden Oasis |
| 3 | 13418779 | 6.620073 | TV,Cable TV,Internet,Wireless Internet,Kitchen... | 4.0 | 1.0 | This light-filled home-away-from-home is super... | 2010-01-01 | 37.772004 | -122.431619 | Beautiful Flat in the Heart of SF! |
| 4 | 3808709 | 4.744932 | TV,Internet,Wireless Internet,Air conditioning... | 2.0 | 1.0 | Cool, cozy, and comfortable studio located in ... | 2010-01-01 | 38.925627 | -77.034596 | Great studio in midtown DC |

5 rows × 118 columns

◀ ━━━━━━━━━ ▶

**Reasoning**: Drop the original columns that were used for creating new features or dummy variables, and other columns not directly relevant for the model.

```
columns_to_drop_original = ['amenities', 'description', 'name', 'host_since', 'id', 'zipcode', 'latitude', 'longitude']
df = df.drop(columns=columns_to_drop_original)
display(df.head())
```

| | log_price | accommodates | bathrooms | number_of_reviews | bedrooms | beds | host_age | Wireless Internet | Kitchen | Heating | ... | host_identity_ve |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 5.010635 | 3.0 | 1.0 | 2.0 | 1.0 | 1.0 | 5632 | 1 | 1 | 1 | ... | |
| 1 | 5.129899 | 7.0 | 1.0 | 6.0 | 3.0 | 3.0 | 5632 | 1 | 1 | 1 | ... | |
| 2 | 4.976734 | 5.0 | 1.0 | 10.0 | 1.0 | 3.0 | 5632 | 1 | 1 | 1 | ... | |
| 3 | 6.620073 | 4.0 | 1.0 | 0.0 | 2.0 | 2.0 | 5632 | 1 | 1 | 1 | ... | |
| 4 | 4.744932 | 2.0 | 1.0 | 4.0 | 0.0 | 1.0 | 5632 | 1 | 1 | 1 | ... | |

5 rows × 110 columns

◀ ━━━━━━━━━ ▶

## Develop model

### Subtask:

Build a regression model to predict listing prices.

**Reasoning**: Define features (X) and target (y), split the data into training and testing sets, import LinearRegression, instantiate the model, and train it.
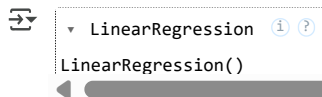
```
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

X = df.drop('log_price', axis=1)
y = df['log_price']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

model = LinearRegression()

model.fit(X_train, y_train)
```

```
▼ LinearRegression  ⓘ ⍰
LinearRegression()
```

## ⌄ Evaluate model

### Subtask:

Evaluate the model's performance using appropriate metrics.

**Reasoning**: Import the necessary metrics for model evaluation.

```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score


import numpy as np
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score

y_pred = model.predict(X_test)

rmse = np.sqrt(mean_squared_error(y_test, y_pred))
mae = mean_absolute_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

print(f"RMSE: {rmse}")
print(f"MAE: {mae}")
print(f"R-squared: {r2}")
```

```
RMSE: 0.4604755088250774
MAE: 0.3496196963521966
R-squared: 0.5808714785079553
```

## ⌄ Summary:

### Data Analysis Key Findings

- The dataset contains 13,251 rows and 29 columns with a mix of numerical and categorical data types.
- Several columns (`host_response_rate`, `review_scores_rating`, `first_review`, `last_review`, `thumbnail_url`, `neighbourhood`) have a high percentage of missing values (ranging from 9.1% to 24.8%).
- Categorical features like `amenities`, `description`, `name`, `neighbourhood`, and `zipcode` exhibit high cardinality.
- Numerical columns like `accommodates`, `bathrooms`, `number_of_reviews`, and `beds` show potential outliers based on their maximum values compared to the 75th percentile.
- Missing values were handled by dropping columns with a high percentage of NaNs and imputing the rest using median for numerical and mode for categorical columns.
- A new feature `host_age` was successfully created by converting the `host_since` column to datetime and calculating the difference in days from a reference date.
- Dummy variables were generated for the 50 most common amenities and other selected categorical features.
- Original columns used for feature engineering or deemed irrelevant (`amenities`, `description`, `name`, `host_since`, `id`, `zipcode`, `latitude`, `longitude`) were dropped.
- A Linear Regression model was trained to predict the `log_price`.
- The trained model achieved an RMSE of approximately 0.460, an MAE of approximately 0.350, and an R-squared score of approximately 0.581 on the test set.

### Insights or Next Steps

- The model explains about 58% of the variance in the log price, suggesting there is room for improvement. Exploring more complex models or additional feature engineering could enhance predictive performance.
- Further investigation into the identified outliers in numerical features might be beneficial to determine if they are valid data points or errors, potentially improving model robustness.

```
print("Shape of the DataFrame:", df.shape)
display(df.head())
```