



DSA5104 Project

Lan Rong	A0169501H
Sun Wei	A0244759N
Wang Yuhua	A0251301B
Yan Yucheng	A0251342R

Contents

- **Clean the Data**
- **Data Generation**
- **Visualization**
- **Relational Database**
- **Non-Relational Database**
- **Queries**



Clean the Data

We first dropped columns that are not useful for our analysis.

```
listings.drop(columns=["scrape_id",  
                      "last_scraped",  
                      "bathrooms",  
                      "calendar_updated",  
                      "availability_30",  
                      "availability_60",  
                      "availability_90",  
                      "availability_365",  
                      "has_availability",  
                      "calendar_last_scraped",  
                      "number_of_reviews",  
                      "number_of_reviews_ltm",  
                      "number_of_reviews_l30d",  
                      "first_review",  
                      "last_review",  
                      "license",  
                      "calculated_host_listings_count",  
                      "calculated_host_listings_count_entire_homes",  
                      "calculated_host_listings_count_private_rooms",  
                      "calculated_host_listings_count_shared_rooms",  
                      "reviews_per_month",  
                      "neighbourhood_group_cleansed"], inplace=True); # Drop columns that are not useful for our analysis
```



Clean the Data

We then fill in the rare null values by getting the information from the url included.

```
listings.loc[listings.host_listings_count.isna(), ['host_name']] = "Jakchai"  
listings.loc[listings.host_listings_count.isna(), ['host_since']] = "2014-01-01"  
listings.loc[listings.host_listings_count.isna(), ['host_location']] = "Bangkok, Thailand"  
listings.loc[listings.host_listings_count.isna(), ['host_is_superhost']] = "f"  
listings.loc[listings.host_listings_count.isna(), ['host_thumbnail_url']] = "https://a0.muscd  
listings.loc[listings.host_listings_count.isna(), ['host_picture_url']] = "https://a0.muscd  
listings.loc[listings.host_listings_count.isna(), ['host_total_listings_count']] = 2  
listings.loc[listings.host_listings_count.isna(), ['host_has_profile_pic']] = 't'  
listings.loc[listings.host_listings_count.isna(), ['host_identity_verified']] = 't'  
listings.loc[listings.host_listings_count.isna(), ['host_listings_count']] = 2
```



Clean the Data

For some numerical null values that are not important, we fill in their mean values.

```
for column in ['minimum_minimum_nights', 'maximum_minimum_nights', 'minimum_maximum_nights', 'maximum_maximum_nights', 'minimum_maximum_nights', 'maximum_minimum_nights', 'minimum_nights', 'maximum_nights', 'number_of_reviews', 'reviews_per_room', 'review_scores_location', 'review_scores_cleanliness', 'review_scores_check_in', 'review_scores_communication', 'review_scores_location', 'review_scores_value']:
    mean = listings[column].mean()
    listings[column].fillna(mean, inplace=True)
```



Clean the Data

Besides, we also dropped neighbourhood, since there is a cleansed column called neighbourhood_cleansed.

We also dropped room_type since it is included in the property_type.

We also dropped the price in listing since it is included in calendar.

We did not drop the properties with price abnormally high since some host may have set an abnormal price in the first place.



Data Generation

We used the review data to generate transaction data.

We used the date from the review table as the `end_date`. After generating some random number as `total_days`, we calculated the `start_date` from `end_date`.

We then used the `adjusted_price` from calendar to generate `total_price`.

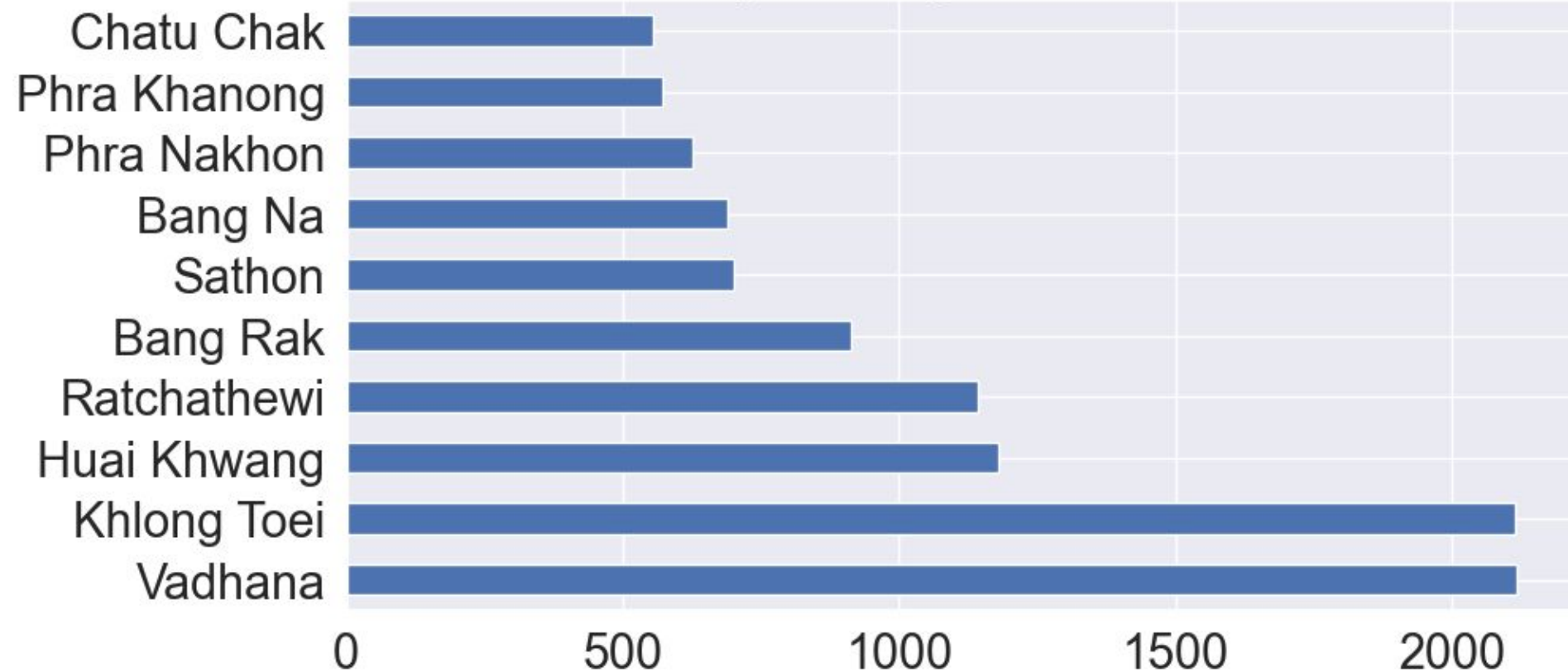
Although all of our transactions include reviews, transactions in real life may not contain a review necessarily.



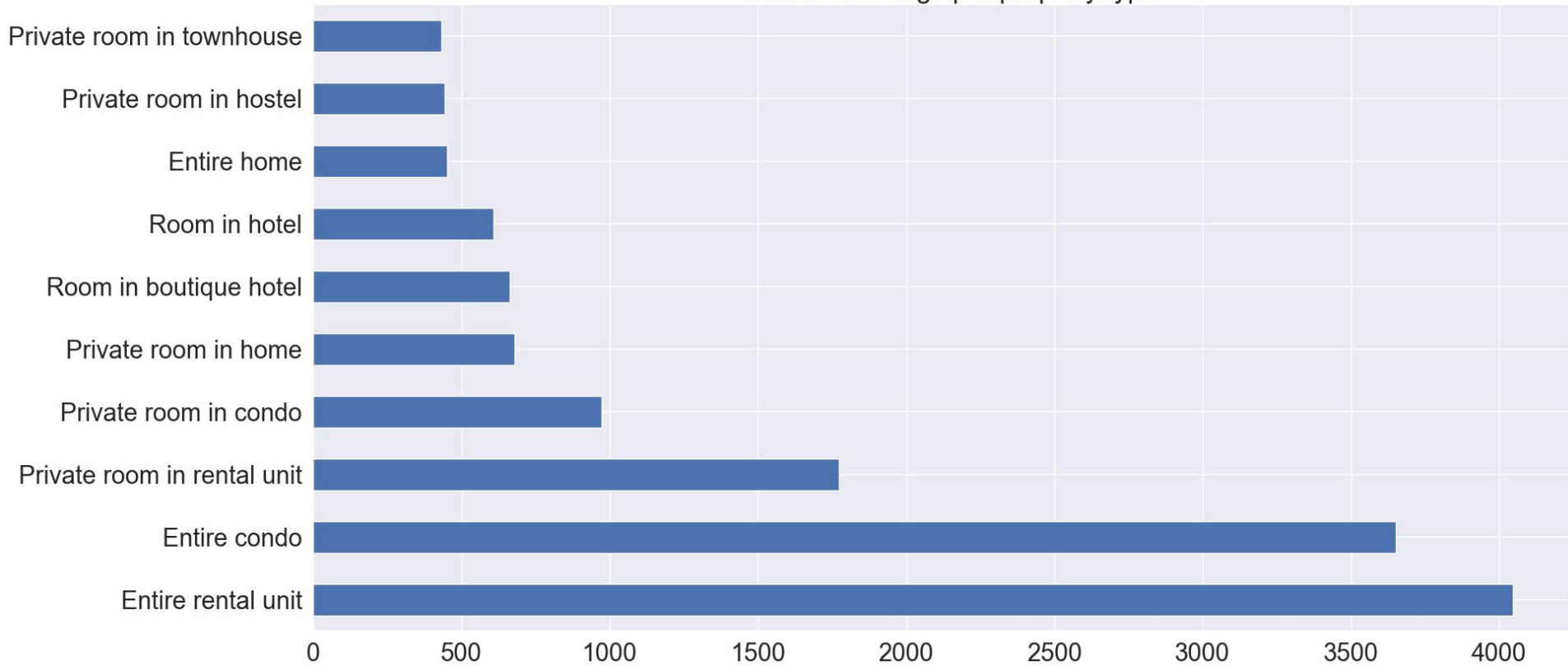
Visualization

We used seaborn to visualize some categorical and numerical features in listing, host, and transaction.

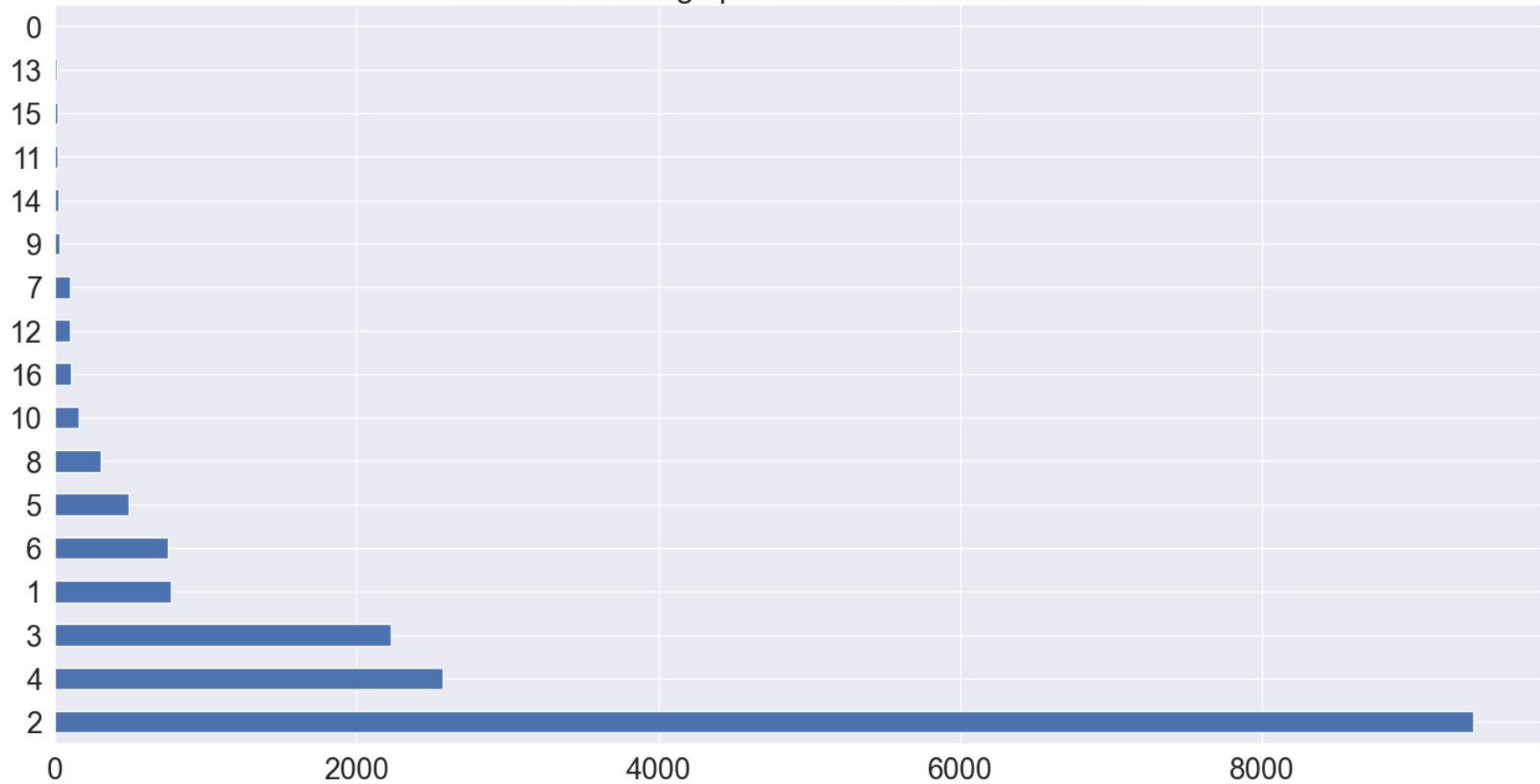
Top 10 neighbourhoods



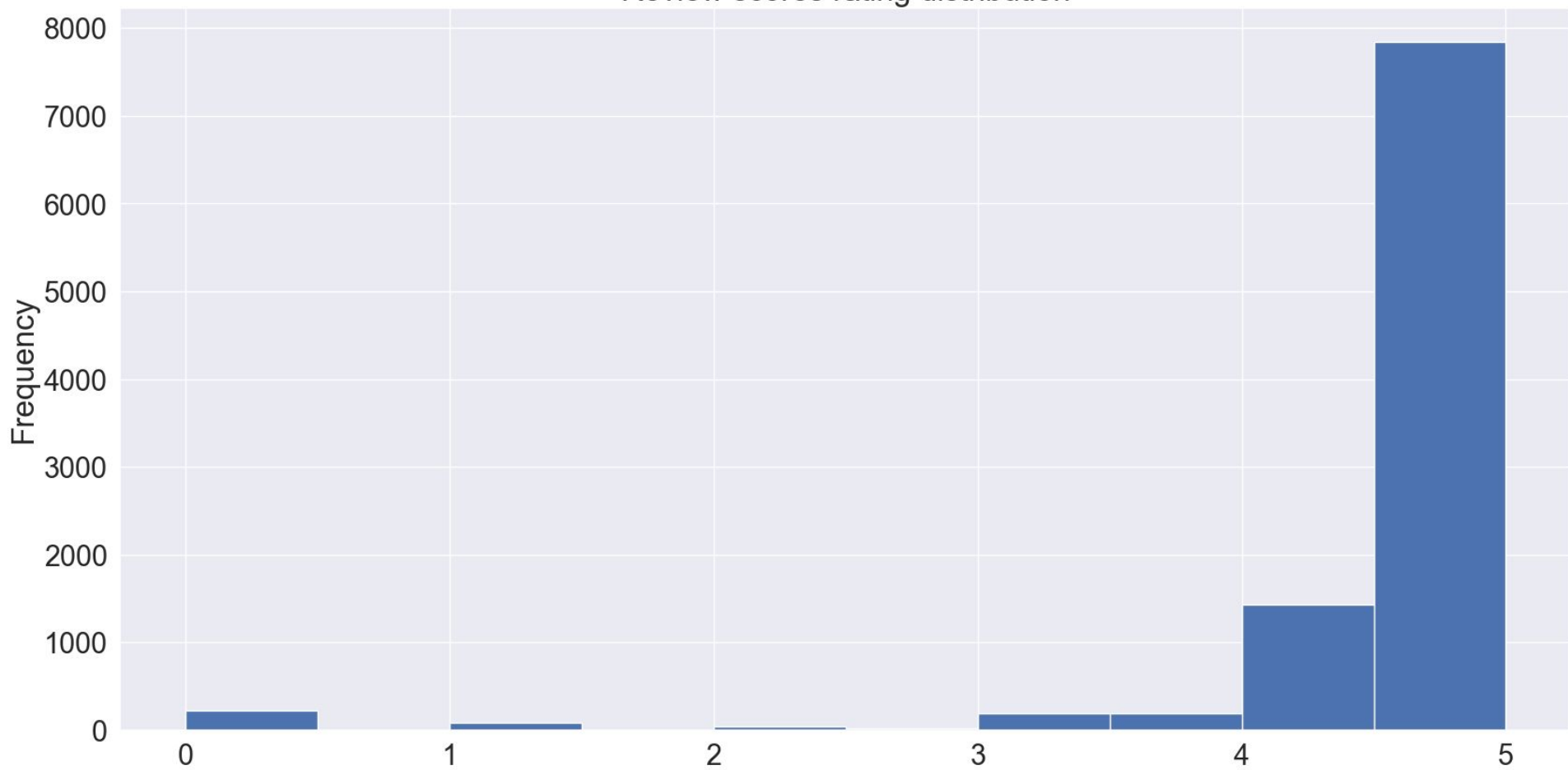
Number of listings per property type



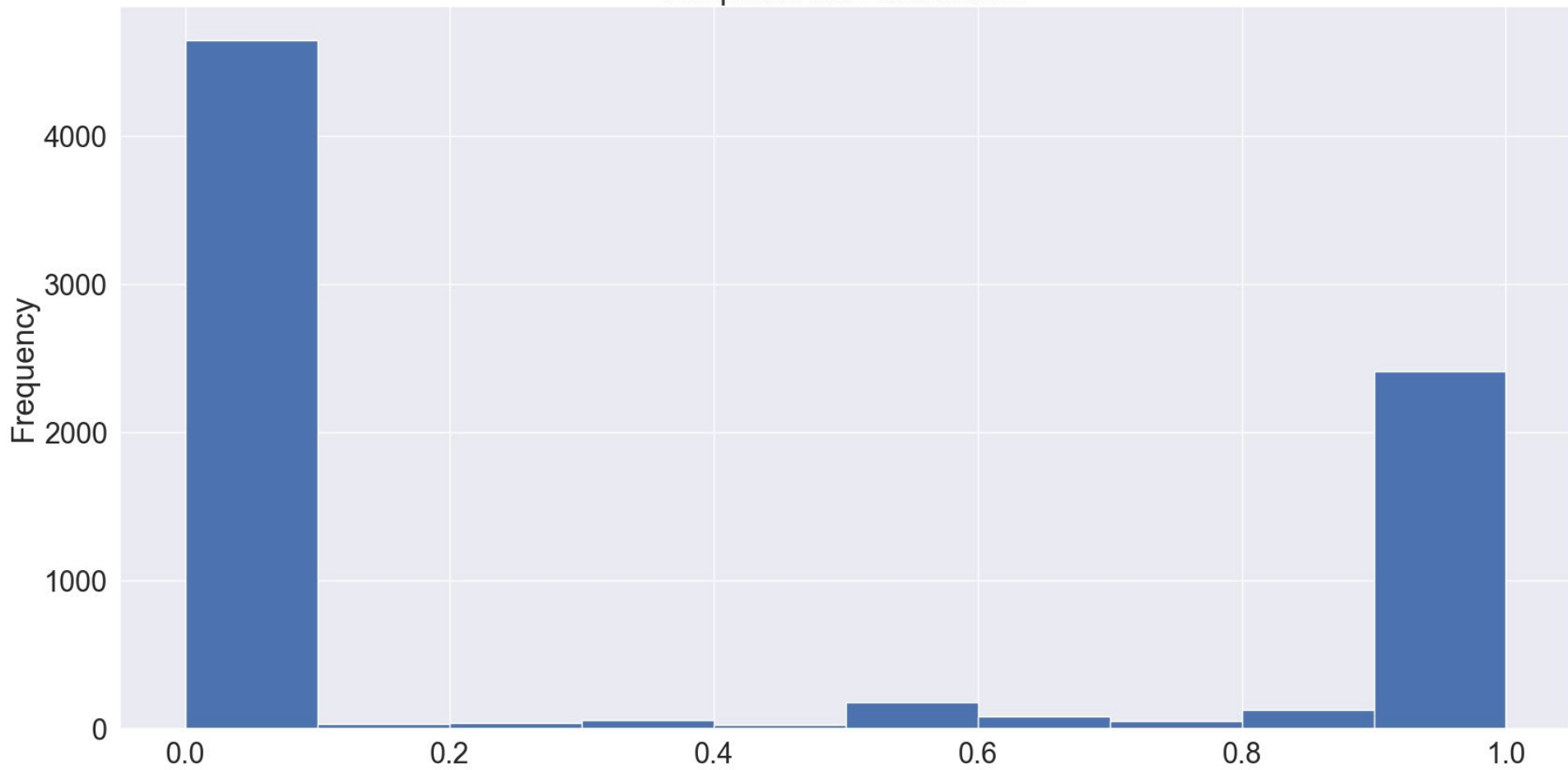
Number of listings per number of accommodates



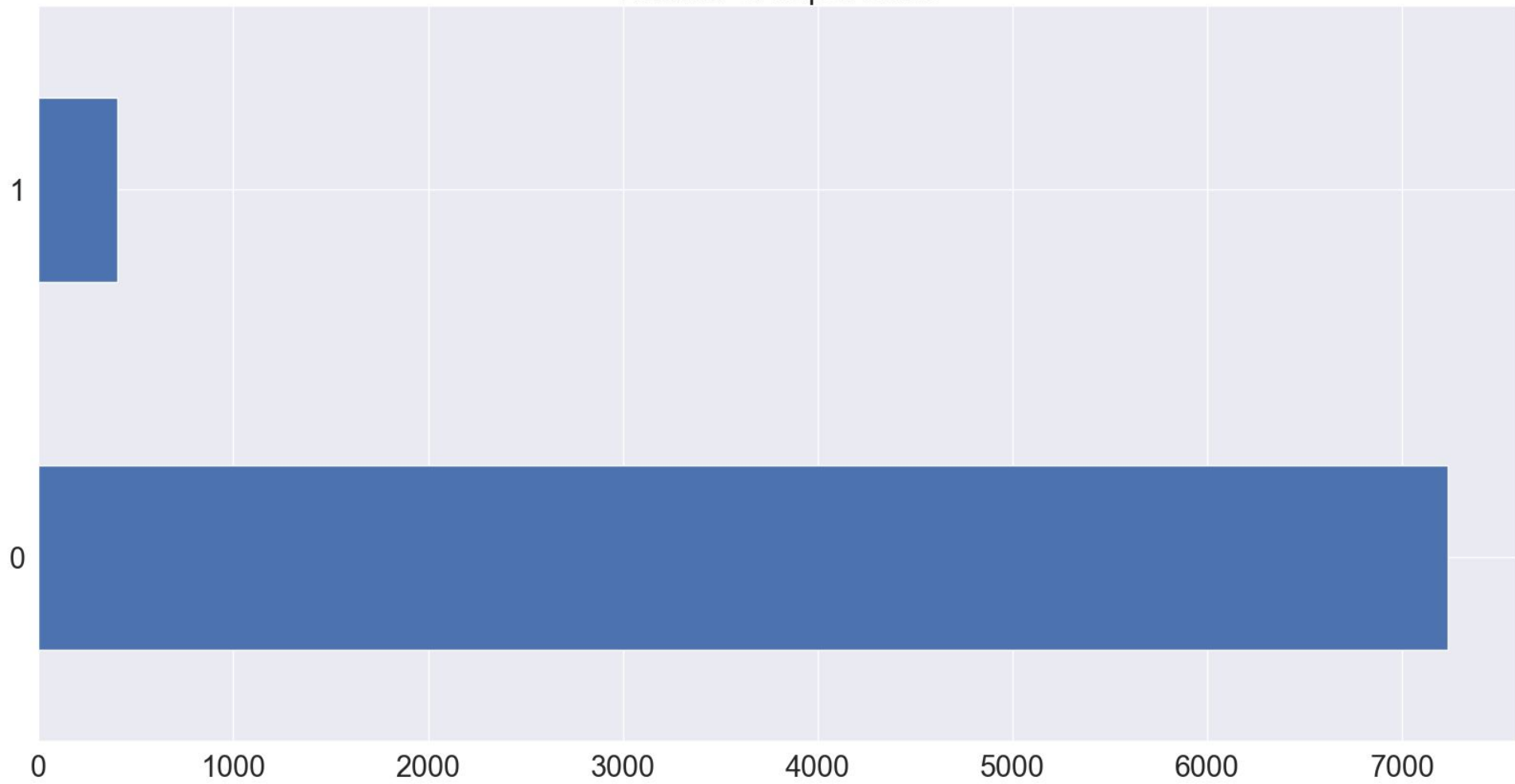
Review scores rating distribution



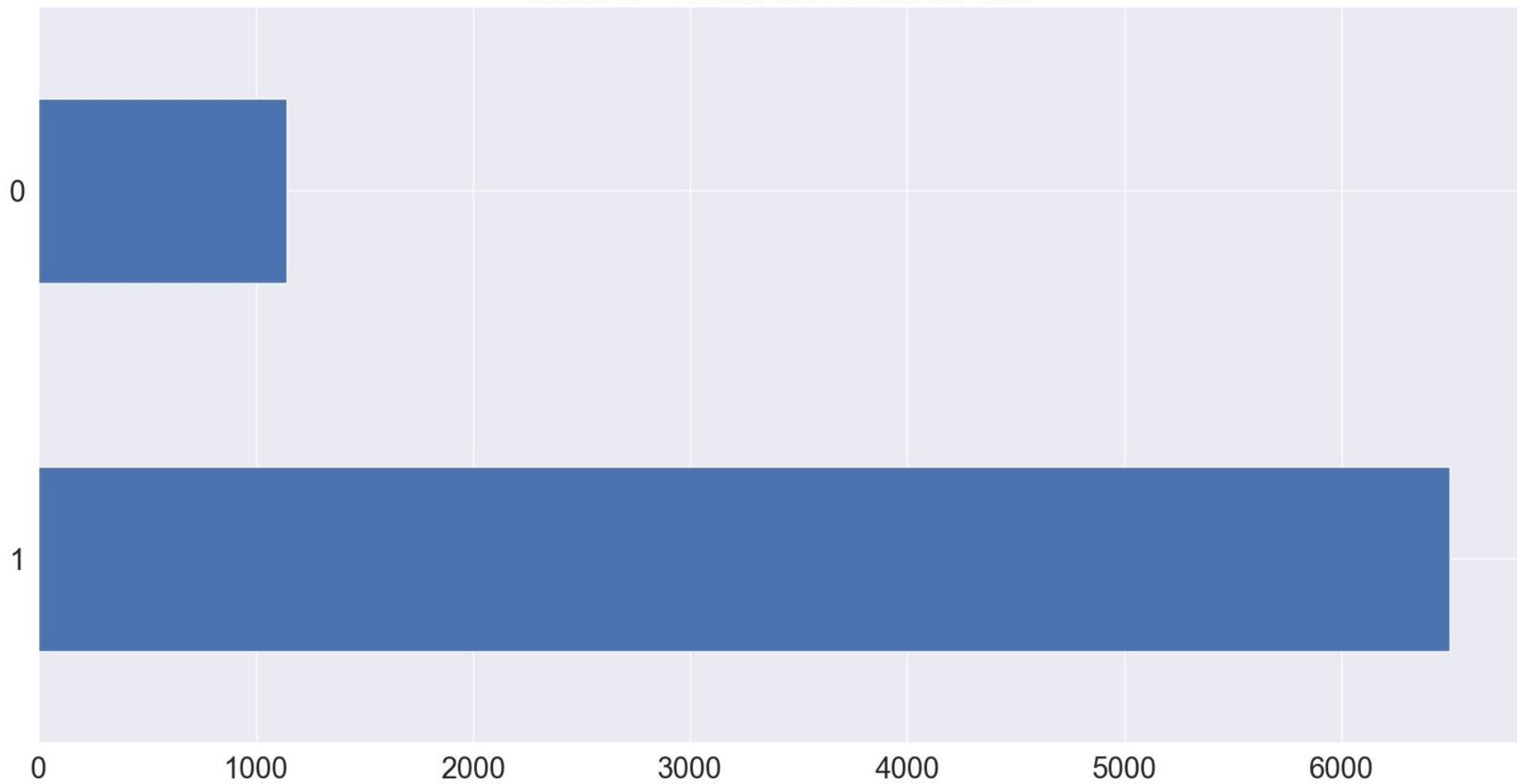
Response rate distribution



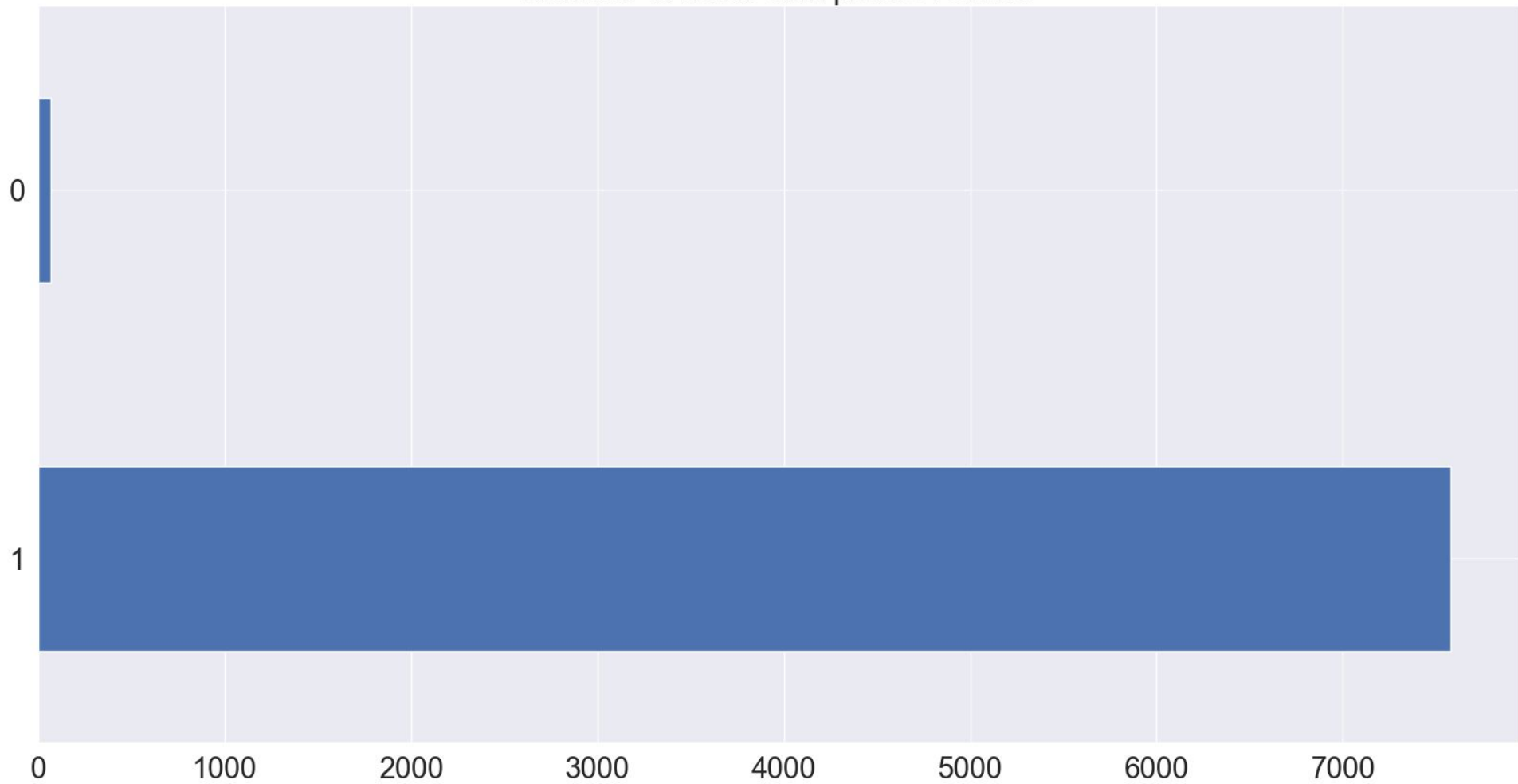
Number of super hosts



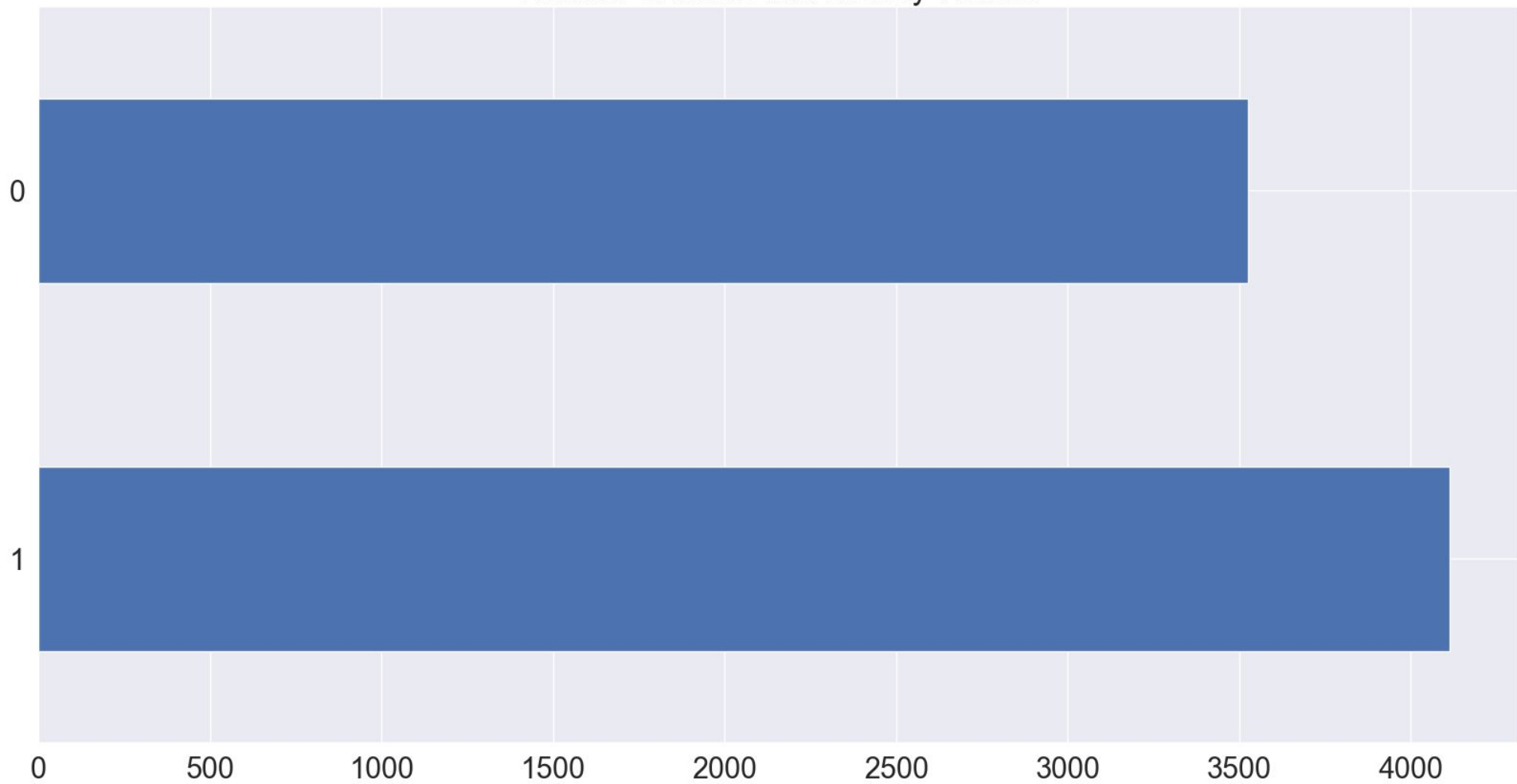
Number of hosts with email verified



Number of hosts with phone verified



Number of hosts with identity verified



Price distribution







Relational Database - Requirement Analysis

- Check property information (availability/comments/...)
- Book available property
- Search property by name/availability/type/area/host/...
- Give comments and reviews after stay
- Search host by ratings/name/...
- Check user information

ER Diagram

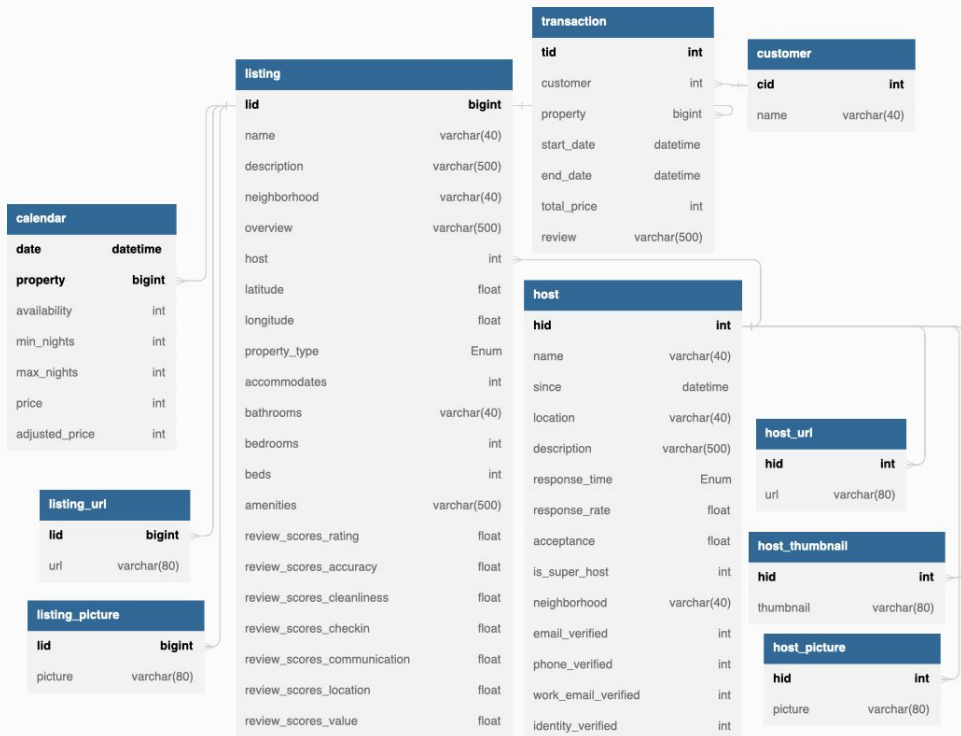


Schema Diagram

When designing schema from ER diagram. We decompose host's and listing's urls and picture urls since they violate the BCNF.

We did not keep any of the relationships since all of them are one-to-many relationships.

You may think that customer, property, and start_date in transaction may violate BCNF. However, one customer may have booked the room and cancelled and booked again in the same day.



Non-Relational Database



For the Non-Relational Database of this project, we chose MongoDB since it is explained in details in the lectures.

We considered various designs for the database. Similar to SQL database, we decide to use the 5 collections because this schema makes it easy to query for the frequently demanded information.

Calendar

Storage size: 136.94 MB
Documents: 6.2 M
Avg. document size: 151.00 B
Indexes: 1
Total index size: 68.64 MB

Customer

Storage size: 4.26 MB
Documents: 210 K
Avg. document size: 30.00 B
Indexes: 1
Total index size: 2.56 MB

Host

Storage size: 1.80 MB
Documents: 76 K
Avg. document size: 716.00 B
Indexes: 1
Total index size: 122.88 kB

Listing

Storage size: 15.64 MB
Documents: 17 K
Avg. document size: 1.96 kB
Indexes: 1
Total index size: 241.66 kB

Transaction

Storage size: 136.87 MB
Documents: 6.2 M
Avg. document size: 151.00 B
Indexes: 1
Total index size: 68.83 MB

MongoDB Schema Design



For Calendars, its arguable that they could be embedded inside the corresponding Listing document. However, when considering the use case in context, we reckon there are mainly 2 scenarios when such information is queried:

1. User wants to see the available dates of a particular listing
2. User wants to find all available listings on a particular date

Embedding Calendars inside Listings made scenarios 1 easy to accomplish. For scenarios 2, on the other hand, we need to iterate through all listings and check if they contain the specified date. This, is computationally expensive. Thus, we decided to store Calendar in its own collection.

Similar use case analysis was performed on Transaction and Listing too.

```
{
  "_id": {
    "$oid": "636b583c69c65a2eb1f7f92b"
  },
  "property": 37433623,
  "date": "2022-08-14",
  "availability": true,
  "min_nights": 7,
  "max_nights": 1125,
  "price": 900,
  "adjusted_price": 900
}
```


Queries

Natural Language	SQL Query	MongoDB Query
Check property information (Find availability of a certain property on a certain day.)	SELECT availability FROM calendar JOIN listing WHERE listing.name=... AND date=...	listing=db.Listing. findOne({"name": ...}) db.Calendar.find({"property": listing._id, "date":...}, {"availability":...})
Book available property (Book a property and set its availability to 0 in booking period)	UPDATE calendar SET availability = 0 WHERE property=... AND date>=... AND date<=...	db.Calendar.upd ateMany({ "property" :listing._id, "date"; {"\$gte": ..., "\$lte":...}}, {"\$set":{"availabili ty":False}})

Queries

Natural Language	SQL Query	MongoDB Query
Search property by name/type/host/...	<pre>SELECT * FROM listing WHERE property_type=... AND neighbourhood=...</pre>	<pre>db.Listing.find({"neighbourhood": ..., "property_type": ..., "bedrooms": ..., "beds": ...})</pre>
Give comments and reviews after stay	<pre>UPDATE transaction SET review=... WHERE tid=...</pre>	<pre>db.Transaction.updateOne ({ "_id":...}, {"\$set": {"review": ...}})</pre>

Queries

Natural Language	SQL Query	MongoDB Query
Search host by ratings/name... (select a host whose response rate is large and his/her listings has high review scores)	<pre>SELECT * FROM host H WHERE response_rate>=... AND EXISTS (SELECT * FROM listing L WHERE L.host=H.hid AND L.review_score>=...)</pre>	<pre>hosts=db.Listing.find({"review_scores_rating": {"\$gte": ...}}, {"host":...}) db.Host.find({"_id": {"\$in": [h.host for h in hosts]}}, "response_rate": {"\$gte:..."}))</pre>
Check user information	<pre>SELECT * FROM customer WHERE name=...</pre>	<pre>db.Customer. find({"name": ...})</pre>



Thank you.

