

Project Team: Walking in Paris

Objective:

Identify the attributes that make Airbnb “Successful” Find the best place to stay in Hawaii

In [1]:

```

1 # Dependencies and Setup
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import numpy as np
5 import requests
6 import time
7
8 # Import statistic library
9 import scipy.stats as stats
10 from scipy.stats import linregress
11
12 # Import seaborn library to create figure
13 import seaborn as sns
14
15 import pprint
16 import json
17 from time import sleep
18 from datetime import date
19 import zipfile
20
21 # Incorporated citipy to determine city based on Latitude and Longitude
22 # pip install citipy
23 from citipy import citipy
24
25 # Output File Location
26 output_folder = "output_data"

```

- 1 Open the dataset from the zip file and ensure memory resources
- 2 Cleaning raw data (airbnb) to delete the column with not values, columns values duplicated, columns with data not required.

```
In [2]: # Open the zip file
2 zf = zipfile.ZipFile('./source_data/detail_listings.zip')
3 df = pd.read_csv(zf.open('detail_listings.csv'), low_memory=False)
4
5 rows_in_df = len(df)
6 print(f'Data set has {rows_in_df} rows')
7
8 # Drop unneeded columns:
9 df = df.drop(columns =
10     "listing_url",
11     "thumbnail_url",
12     "picture_url",
13     "medium_url",
14     "xl_picture_url",
15     "host_url",
16     "host_thumbnail_url",
17     "host_picture_url",
18     "scrape_id",
19     "host_listings_count", # There is a host_total_listings_count that is a duplicate
20     "host_acceptance_rate", # NaN
21     "calendar_last_scraped",
22     "bed_type",
23     "last_scraped",
24     "calendar_updated",
25     "has_availability",
26     "availability_30",
27     "availability_60",
28     "availability_90",
29     "availability_365",
30     "license"
31 )
32
33 df.head(1)
```

Data set has 23745 rows

| | id | name | summary | space | description | experiences_offered | neighborhood_overview | notes | location |
|---|-----------|-------------|---|---|---|----------------------------|---|-------------------|-----------------|
| 0 | 5065 | MAUKA BB | Perfect for your vacation, Staycation or just ... | Mauka B & B is 400 square feet studio with a p... | Perfect for your vacation, Staycation or just ... | none | Neighbors here are friendly but are not really... | Kalo Sta Pe w 100 | Locat near |

1 rows × 85 columns

In [3]:

```
1 # Group the columns and show the list of the columns
2
3 print('')
4 print('-----')
5 print('Property-specific columns:')
6 print('-----')
7 # These are wildcard search terms:
8 cols = ["experiences_offered", "street", "city", "state", "zipcode", "market"
9       , "location", "accomodate", "room", "guest", "is_business_travel_ready",
10      for col in df.columns:
11          if col in cols:
12              print(col)
13
14 print('')
15 print('-----')
16 print('All columns:')
17 print('-----')
18 for col in df.columns:
19     print(col)
```

Property-specific columns:

experiences_offered
street
city
state
zipcode
market
country
instant_bookable
is_business_travel_ready

All columns:

id
name
summary
space
description
experiences_offered
neighborhood_overview
notes
transit
access
interaction
house_rules
host_id
host_name
host_since
host_location
host_about
host_response_time
host_response_rate

host_is_superhost
host_neighbourhood
host_total_listings_count
host_verifications
host_has_profile_pic
host_identity_verified
street
neighbourhood
neighbourhood_cleansed
neighbourhood_group_cleansed
city
state
zipcode
market
smart_location
country_code
country
latitude
longitude
is_location_exact
property_type
room_type
accommodates
bathrooms
bedrooms
beds
amenities
square_feet
price
weekly_price
monthly_price
security_deposit
cleaning_fee
guests_included
extra_people
minimum_nights
maximum_nights
minimum_minimum_nights
maximum_minimum_nights
minimum_maximum_nights
maximum_maximum_nights
minimum_nights_avg_ntm
maximum_nights_avg_ntm
number_of_reviews
number_of_reviews_ltm
first_review
last_review
review_scores_rating
review_scores_accuracy
review_scores_cleanliness
review_scores_checkin
review_scores_communication
review_scores_location
review_scores_value
requires_license
jurisdiction_names
instant_bookable

```
is_business_travel_ready  
cancellation_policy  
require_guest_profile_picture  
require_guest_phone_verification  
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
```

```
In [4]: # Create a dataframe with 1 row for each column in the source:
1 data_dict = pd.DataFrame(columns = [
2     "ColumnNumber"
3     , "Column"
4     , "Group"
5     , "Subgroup"
6     , "Notes"
7 ])
8 ]
9
10 # Loop through
11 for column in df:
12     data_dict = data_dict.append({
13         "Column": column
14         , "Group": ""
15         , "Subgroup": ""
16         , "Notes": ""
17     }
18     , ignore_index=True
19 )
20
21 # Make the index the column number
22 data_dict['ColumnNumber'] = data_dict.index
23
24 # Add context for each column
25 for index, row in data_dict.iterrows():
26     val = str(row["Column"])
27
28     if val == "id":
29         data_dict.loc[index, "Notes"] = "Primary Key"
30     # Val.startswith("host_")
31     if val.find("host") > -1:
32         data_dict.loc[index, "Group"] = "Host-related"
33
34     if val.find("review") > -1:
35         data_dict.loc[index, "Group"] = "Review-related"
36
37     if (val.find("rule") > -1) | (val.find("require") > -1):
38         data_dict.loc[index, "Group"] = "Requirements-related"
39
40     if (val.find("neighbo") > -1):
41         data_dict.loc[index, "Group"] = "Neighborhood-related"
42
43     if (val.find("nights") > -1):
44         data_dict.loc[index, "Group"] = "Booking-related"
45
46 # Property specific:
47 cols = [
48     "property", "room", "guest", "is_business_travel_ready", "instant_bookable",
49     , "space", "name", "summary", "description", "notes", "transit", "accessories",
50     , "property_type", "room_type", "amenities"
51 ]
52     if any(word in val for word in cols):
53         data_dict.loc[index, "Group"] = "Property-related"
54
55 # Property - "Location" specific:
56 cols = [
```

```

57     "location", "street", "city", "state", "zipcode", "market", "country"
58     , "latitude", "longitude", "is_location_exact"
59 ]
60 if any(word in val for word in cols):
61     data_dict.loc[index, "Group"] = "Property-related"
62     data_dict.loc[index, "Subgroup"] = "Location"
63
64 # Property - "Interior" specific:
65 cols = [
66     "bathrooms", "bedrooms", "beds", "square_feet"
67 ]
68 if any(word in val for word in cols):
69     data_dict.loc[index, "Group"] = "Property-related"
70     data_dict.loc[index, "Subgroup"] = "Interior"
71
72 # Property - "Pricing" specific:
73 # Things about price, or that we might associate with price
74 cols = [
75     "price", "weekly_price", "monthly_price"
76     , "security_deposit", "cleaning_fee"
77     , "accommodates", "guests_included", "extra_people"
78     , "experiences_offered"
79 ]
80 if any(word in val for word in cols):
81     data_dict.loc[index, "Group"] = "Property-related"
82     data_dict.loc[index, "Subgroup"] = "Pricing"
83
84 # Add additional code to accomodate a group and subgroup of properties:
85 if val == "accommodates":
86     data_dict.loc[index, "Group"] = "Property-related"
87     data_dict.loc[index, "Subgroup"] = "Pricing"
88
89 print(data_dict.groupby(["Group", "Subgroup"])["Column"].count())
90 print('')
91 print('Columns without a group:')
92 data_dict[data_dict["Group"] == ""]
93

```

| Group | Subgroup | |
|----------------------------|----------|----|
| | | 1 |
| Booking-related | | 8 |
| Host-related | | 12 |
| Neighborhood-related | | 5 |
| Property-related | | 20 |
| | Interior | 4 |
| | Location | 13 |
| | Pricing | 9 |
| Requirements-related | | 2 |
| Review-related | | 11 |
| Name: Column, dtype: int64 | | |

Columns without a group:

| Out[4]: | ColumnNumber | Column | Group | Subgroup | Notes |
|---------|--------------|--------|-------|----------|-------------|
| | 0 | 0 | id | | Primary Key |

Making additional cleaning of the raw data in the rows

- Drop rows with number_of_reviews = 0
- Get names of indexes for which column number_of_reviews has value 0
- Strip out the dollar symbol (\$) from the pricing based columns:
price", "weekly_price", "monthly_price", "security_deposit", "cleaning_fee", "extra_people"

In [5]:

```
1 # Drop by index:
2 indexNames = df[df['number_of_reviews'] == 0 ].index
3 df.drop(indexNames , inplace=True)
4
5 # Drop rows with number_of_reviews = NaN
6 df.dropna(axis=0, subset=('number_of_reviews', ))
7
8 # Drop rows with NaN review data
9 df.dropna(axis=0, subset=(
10     'review_scores_rating',
11     'review_scores_accuracy',
12     'review_scores_cleanliness',
13     'review_scores_checkin',
14     'review_scores_communication',
15     'review_scores_location',
16     'review_scores_value'
17 ))
18
19 print(f'Remaining rows after cleanup: {len(df)}')
```

Remaining rows after cleanup: 18291

```
In [6]: # Strip the $ out of the pricing-based columns
1 data_dict[(data_dict["Group"] == "Property-related") & (data_dict["Subgroup"]
2
3
4 if df["price"].dtype != "float64":
5     df["price"] = df["price"].str.replace("$", "")
6     df["price"] = df["price"].str.replace(",", "")
7     df["price"] = df["price"].astype("float64")
8
9 if df["weekly_price"].dtype != "float64":
10    df["weekly_price"] = df["weekly_price"].str.replace("$", "")
11    df["weekly_price"] = df["weekly_price"].str.replace(",", "")
12    df["weekly_price"] = df["weekly_price"].astype("float64")
13
14 if df["monthly_price"].dtype != "float64":
15    df["monthly_price"] = df["monthly_price"].str.replace("$", "")
16    df["monthly_price"] = df["monthly_price"].str.replace(",", "")
17    df["monthly_price"] = df["monthly_price"].astype("float64")
18
19 if df["security_deposit"].dtype != "float64":
20    df["security_deposit"] = df["security_deposit"].str.replace("$", "")
21    df["security_deposit"] = df["security_deposit"].str.replace(",", "")
22    df["security_deposit"] = df["security_deposit"].astype("float64")
23
24 if df["cleaning_fee"].dtype != "float64":
25    df["cleaning_fee"] = df["cleaning_fee"].str.replace("$", "")
26    df["cleaning_fee"] = df["cleaning_fee"].str.replace(",", "")
27    df["cleaning_fee"] = df["cleaning_fee"].astype("float64")
28
29 if df["extra_people"].dtype != "float64":
30    df["extra_people"] = df["extra_people"].str.replace("$", "")
31    df["extra_people"] = df["extra_people"].str.replace(",", "")
32    df["extra_people"] = df["extra_people"].astype("float64")
33
34 data_dict
```

| ColumnNumber | Column | Group | Subgroup | Notes |
|--------------|--|------------------|----------|-------------|
| 0 | id | | | Primary Key |
| 1 | name | Property-related | | |
| 2 | summary | Property-related | | |
| 3 | space | Property-related | | |
| 4 | description | Property-related | | |
| ... | ... | ... | ... | ... |
| 80 | calculated_host_listings_count | Host-related | | |
| 81 | calculated_host_listings_count_entire_homes | Host-related | | |
| 82 | calculated_host_listings_count_private_rooms | Property-related | | |

| ColumnNumber | Column | Group | Subgroup | Notes |
|--------------|---|------------------|----------|-------|
| 83 | calculated_host_listings_count_shared_rooms | Property-related | | |
| 84 | reviews_per_month | Review-related | | |

85 rows × 5 columns

```

1 # Create the bins to separate the number of properties by host
2
3 bins = [1, 60, 119, 179, 239, 299]
4
5 # labels - always one less than the range
6 group_names = ["0-59 properties", "60-119 properties", "120-179 properties",
7 "180-239 properties", "240+ properties"]
8 df["Bin_NumProperties"] = pd.cut(df["host_total_listings_count"], bins,
9 labels=group_names)
10
11 df_analyze = df[["host_id", "Bin_NumProperties", "host_total_listings_count"]]
12 df_analyze = df_analyze.rename(columns={"host_total_listings_count":
13 "listings"})
14
15 # Loop through and get aggs:
16 for index, row in df_analyze.iterrows():
17     df_analyze.loc[index, 'reviews'] = df_analyze["number_of_reviews"].sum()
18
19 #df_analyze["reviews"] =
20 df_analyze["rating_mean"] = df_analyze["review_scores_rating"].mean()
21 df_analyze["rating_value_mean"] = df_analyze["review_scores_value"].mean()
22 df_analyze.head()

```

In [7]: 1

```
In [8]: # Create the bins to separate the number of properties by host
1 bins = [0, 50, 100, 150, 200, 250, 299]
2
3
4 # Labels - always one less than the range
5 group_names = ["0-50 properties", "51-100 properties", "101-150 properties",
6 "201-250 properties", "251+ properties"]
7
8 df["Bin_NumProperties"] = pd.cut(df["calculated_host_listings_count"], bins,
9
10
11 # Loop through and get aggs:
12 total_reviews_per_host = df.groupby("host_id").sum()["number_of_reviews"]
13 average_review_per_host = df.groupby("host_id").mean()["review_scores_rating"]
14 rating_mean_by_bin = df.groupby("Bin_NumProperties").mean()["review_scores_ra
15
16 #df_analyze["rating_value_mean"] = df_analyze["review_scores_value"].mean()
17 df.head(1)
```

Out[8]:

| | id | name | summary | space | description | experiences_offered | neighborhood_overview | not |
|---|-----------|-------------|---|---|---|----------------------------|---|-------------------------------|
| 0 | 5065 | MAUKA BB | Perfect for your vacation, Staycation or just ... | Mauka B & B is 400 square feet studio with a p... | Perfect for your vacation, Staycation or just ... | none | Neighbors here are friendly but are not really... | Locati near Kalo Sta Pe w 100 |

1 rows × 86 columns

| | |
|---|--|
| 1 | |
|---|--|

```
In [10]: # data_dict[data_dict["Group"] == "Host-related"]
1 host_group = df.groupby("neighbourhood").agg({
2     "number_of_reviews": [
3         np.count_nonzero,
4         np.mean,
5         np.median,
6         np.var,
7         np.std
8     ]
9 })
10 )
11 host_group.rename(columns={'count_nonzero': 'Reviews'})
12 host_group.head()
13 # bin on reviews_per_month
```

Out[10]:

| | | number_of_reviews | | | | |
|--|-----------------------------------|-------------------|-----------|--------|--------------|------------|
| | | count_nonzero | mean | median | var | std |
| | neighbourhood | | | | | |
| | Aiea | 14 | 40.428571 | 8.0 | 2872.417582 | 53.594940 |
| | Airport | 3 | 4.000000 | 4.0 | 9.000000 | 3.000000 |
| | Ala Moana/Kakaako | 28 | 23.214286 | 8.5 | 1138.767196 | 33.745625 |
| | Aliamanu/Salt Lake/Foster Village | 8 | 74.875000 | 9.5 | 17739.267857 | 133.188843 |
| | Central Oahu | 43 | 29.302326 | 17.0 | 1214.073090 | 34.843552 |

```
In [11]: 1 bins_num_reviews = [10, 20, 50, 100, 200, 300, 350, 400, 450, 500, 550, 600,
2
3 # Labels - always one less than the range
4 group_names_num_reviews = [
5     "1-10"
6     , "11-20"
7     , "21-50"
8     , "100-200"
9     , "201-250"
10    , "251-300"
11    , "301-350"
12    , "351-400"
13    , "401-550"
14    , "451-500"
15    , "501-550"
16    , "551-600"
17 ]
18
19 df["Bin_NumReviews"] = pd.cut(
20     df["number_of_reviews"],
21     bins = bins_num_reviews,
22     labels = group_names_num_reviews
23 )
24
25 df["Bin_NumReviews"].value_counts()
26
```

```
Out[11]: 11-20      3790
1-10       2526
21-50      2368
100-200    1235
201-250    250
251-300    32
301-350    23
351-400    10
401-550    6
551-600    3
451-500    3
501-550    1
Name: Bin_NumReviews, dtype: int64
```

```
In [12]: 1 print(f'Sum: {df["number_of_reviews"].sum()}')
2 print(f'Count: {df["number_of_reviews"].count()}')
3
4 df.groupby("Bin_NumReviews")["number_of_reviews"].sum()
```

Sum: 623494

Count: 18291

```
Out[12]: Bin_NumReviews
1-10          37991
11-20         126389
21-50         169320
100-200        169224
201-250        59598
251-300        10221
301-350        8452
351-400        4217
401-550        2864
451-500        1561
501-550        569
551-600        1983
Name: number_of_reviews, dtype: int64
```

```
In [13]: df.groupby("Bin_NumReviews").agg({
    "number_of_reviews": [
        np.count_nonzero,
        np.mean,
        np.median,
        np.var,
        np.std
    ]
})
```

Out[13]:

| Bin_NumReviews | number_of_reviews | | | | | |
|----------------|-------------------|------------|--------|------------|-----------|--|
| | count_nonzero | mean | median | var | std | |
| 1-10 | 2526 | 15.039984 | 15.0 | 8.161173 | 2.856777 | |
| 11-20 | 3790 | 33.348021 | 32.0 | 74.548947 | 8.634173 | |
| 21-50 | 2368 | 71.503378 | 70.0 | 197.305439 | 14.046545 | |
| 100-200 | 1235 | 137.023482 | 131.0 | 741.095883 | 27.223076 | |
| 201-250 | 250 | 238.392000 | 232.0 | 802.456161 | 28.327657 | |
| 251-300 | 32 | 319.406250 | 318.0 | 176.958669 | 13.302581 | |
| 301-350 | 23 | 367.478261 | 366.0 | 157.988142 | 12.569333 | |
| 351-400 | 10 | 421.700000 | 421.5 | 252.900000 | 15.902830 | |
| 401-550 | 6 | 477.333333 | 485.5 | 246.666667 | 15.705625 | |
| 451-500 | 3 | 520.333333 | 518.0 | 24.333333 | 4.932883 | |
| 501-550 | 1 | 569.000000 | 569.0 | NaN | NaN | |
| 551-600 | 3 | 661.000000 | 669.0 | 577.000000 | 24.020824 | |

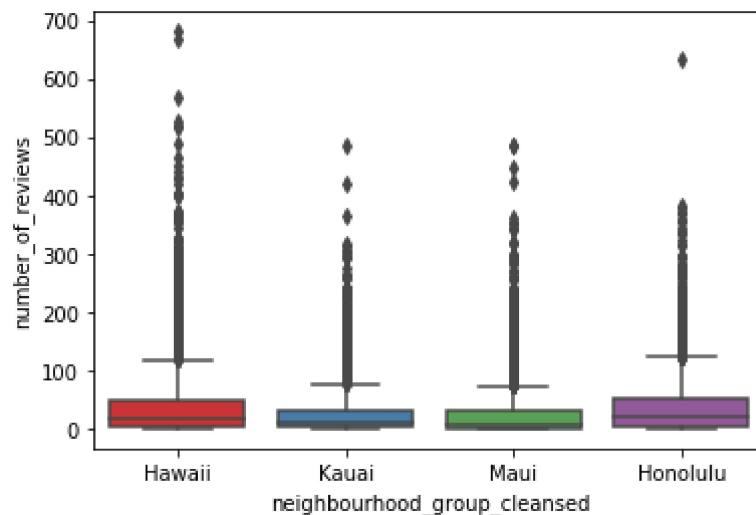
```
1 Hypothesis candidates:
2 Neighborhood
3     - Is rating related to price within a neighborhood?
4     - The highest price properties in a neighborhood will have lower ratings
5     - Best neighborhoods by rating
6     - Best neighborhoods by price
7     - Listing types by neighborhood
8     - Listing types by price
9     - Listing types by rating
10 By island:
11     - Is one island more expensive than the others?
```

```
In [ ]: 1 START QUESTIONS CODE PER TEAM MEMBER
```

```
1 Scott:
2
```

```
In [14]: # Use seaborn Library to create a boxplot
# sns.boxplot(x="day", y="total_bill", hue="smoker", data=df, palette="Set1")
sns.boxplot(
    x = "neighbourhood_group_cleansed",
    y = "number_of_reviews",
    #hue = "price",
    data = df,
    palette = "Set1"
)
# Bin by island
# data_dict[(data_dict["Group"] == "Property-related") & (data_dict["Subgroup"] == "Honolulu")]
```

Out[14]: <matplotlib.axes._subplots.AxesSubplot at 0x2519f4722c8>



```
In [15]: from pandas.api.types import is_string_dtype
from pandas.api.types import is_numeric_dtype
is_string_dtype(df['price'])
is_numeric_dtype(df['price'])
df["price"]
```

Out[15]: 0 85.0
1 110.0
2 92.0
3 299.0
4 92.0
...
23569 111.0
23576 79.0
23582 70.0
23656 146.0
23741 499.0
Name: price, Length: 18291, dtype: float64

In [16]:

```

1 # Bin by island
2 isle = df[['latitude', 'longitude']].head(50).copy()
3 isle["Island"] = ""
4
5 curr_rows = len(isle)
6 print(f'There are {curr_rows} coordinates in the dataframe')
7
8 df["neighbourhood_group_cleansed"]

```

There are 50 coordinates in the dataframe

Out[16]:

```

0 Hawaii
1 Hawaii
2 Hawaii
3 Hawaii
4 Kauai
...
23569 Maui
23576 Hawaii
23582 Honolulu
23656 Honolulu
23741 Kauai
Name: neighbourhood_group_cleansed, Length: 18291, dtype: object

```

Below we are testing the hypothesis 1: If a host has a high number of properties then their ratings are not higher than hosts with fewer properties

| | |
|---|--|
| 1 | Meg Section |
| 2 | Evaluate the relationship between "age of property" and review_scores_rating |
| 3 | Evaluate the relationship between "price" and review_scores_rating |
| 4 | Evaluate the relationship between "how many properties someone has" and review_scores_rating |
| 5 | Evaluate the relationship between "number of bedrooms" and review_scores_rating |
| 6 | |
| 7 | Observations: |

In [17]:

```

1 # Define the bins by number of properties by host
2 bins = [0, 50, 100, 150, 200, 250, 299]
3
4 # Labels - always one less than the range
5 group_names = ["0-50 properties", "51-100 properties", "101-150 properties",
6                 "201-250 properties", "251+ properties"]
7
8 df["Bin_NumProperties"] = pd.cut(df["calculated_host_listings_count"], bins,
9
10 rating_mean_by_bin = df.groupby("Bin_NumProperties").mean()[["review_scores_ra
11
12 df.head(2)

```

Out[17]:

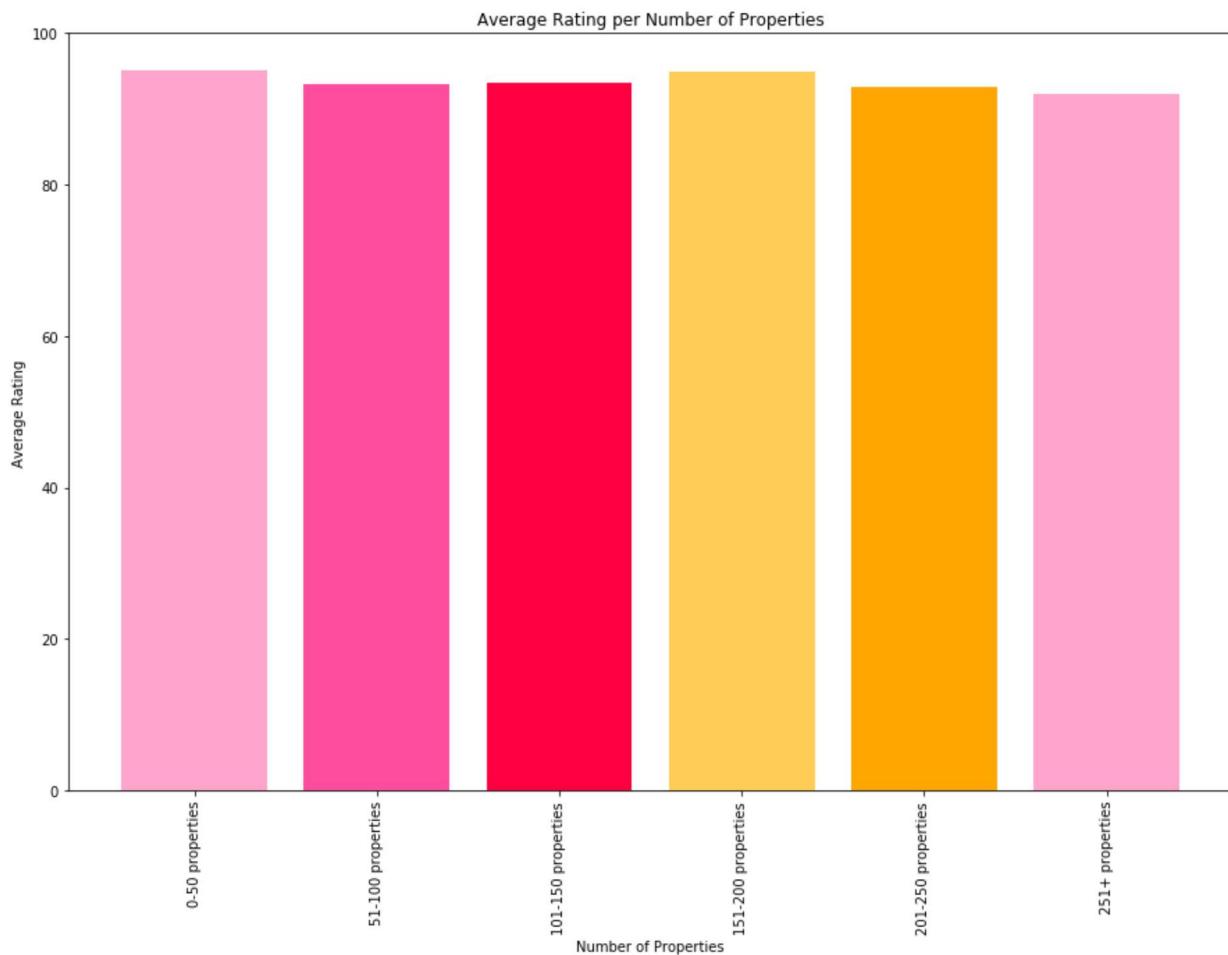
| | | id | name | summary | space | description | experiences_offered | neighborhood_overview |
|---|------|---|---|---|---|-------------|---|-----------------------|
| 0 | 5065 | MAUKA BB | Perfect for your vacation, Staycation or just ... | Mauka B & B is 400 square feet studio with a p... | Perfect for your vacation, Staycation or just ... | none | Neighbors here are friendly but are not really... | |
| 1 | 5269 | Upcountry Hospitality in the 'Auwai Suite | The 'Auwai Suite is a lovely, self-contained a... | The 'Auwai Suite is a lovely, self-contained a... | The 'Auwai Suite is a lovely, self-contained a... | none | We are located on the "sunny side" of Waimea, ... | |

2 rows × 87 columns

In [18]:

```
1 # Define the axis
2 x_axis = group_names
3 y_axis = rating_mean_by_bin
4
5 # Bring in the colors from Hawaii
6 hawaiian_hibiscus_hex_palette = ["#ffa4cd", "#ff4e9e", "#ff0040", "#ffcd55",
7
8 # Define the properties of the figure
9
10 plt.subplots(figsize=(15,10))
11 plt.title("Average Rating per Number of Properties")
12 plt.xlabel("Number of Properties")
13 plt.ylabel("Average Rating")
14 plt.xticks(rotation=90)
15 plt.ylim(0, 100)
16 plt.bar(x_axis, y_axis, color = hawaiian_hibiscus_hex_palette)
17
18 # Show the figure
19 plt.show
```

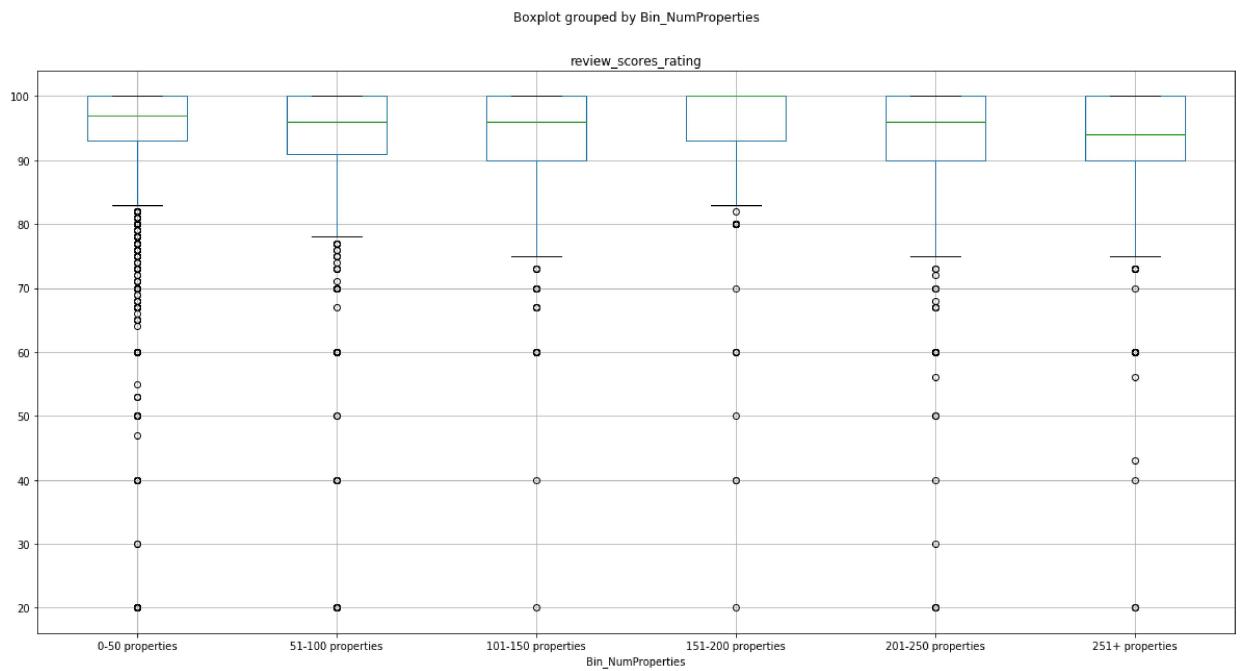
Out[18]: <function matplotlib.pyplot.show(*args, **kw)>



```
1 #ask Meg the objective of the figure below
```

In [19]: 1 df.boxplot("review_scores_rating", by="Bin_NumProperties", figsize=(20, 10))

Out[19]: <matplotlib.axes._subplots.AxesSubplot at 0x251a2053788>



To run the next code, we have "import scipy.stats as stats". Next we are creating a set of group by Number of properties and review score rating

In [20]:

```
1 # Extract individual groups
2 group1 = df[df["Bin_NumProperties"] == "0-50 properties"]["review_scores_ratio"]
3 group2 = df[df["Bin_NumProperties"] == "51-100 properties"]["review_scores_ratio"]
4 group3 = df[df["Bin_NumProperties"] == "101-150 properties"]["review_scores_ratio"]
5 group4 = df[df["Bin_NumProperties"] == "151-200 properties"]["review_scores_ratio"]
6 group5 = df[df["Bin_NumProperties"] == "201-250 properties"]["review_scores_ratio"]
7 group6 = df[df["Bin_NumProperties"] == "251+ properties"]["review_scores_ratio"]
8
9 stats.f_oneway(group1, group2, group3, group4, group5, group6)
```

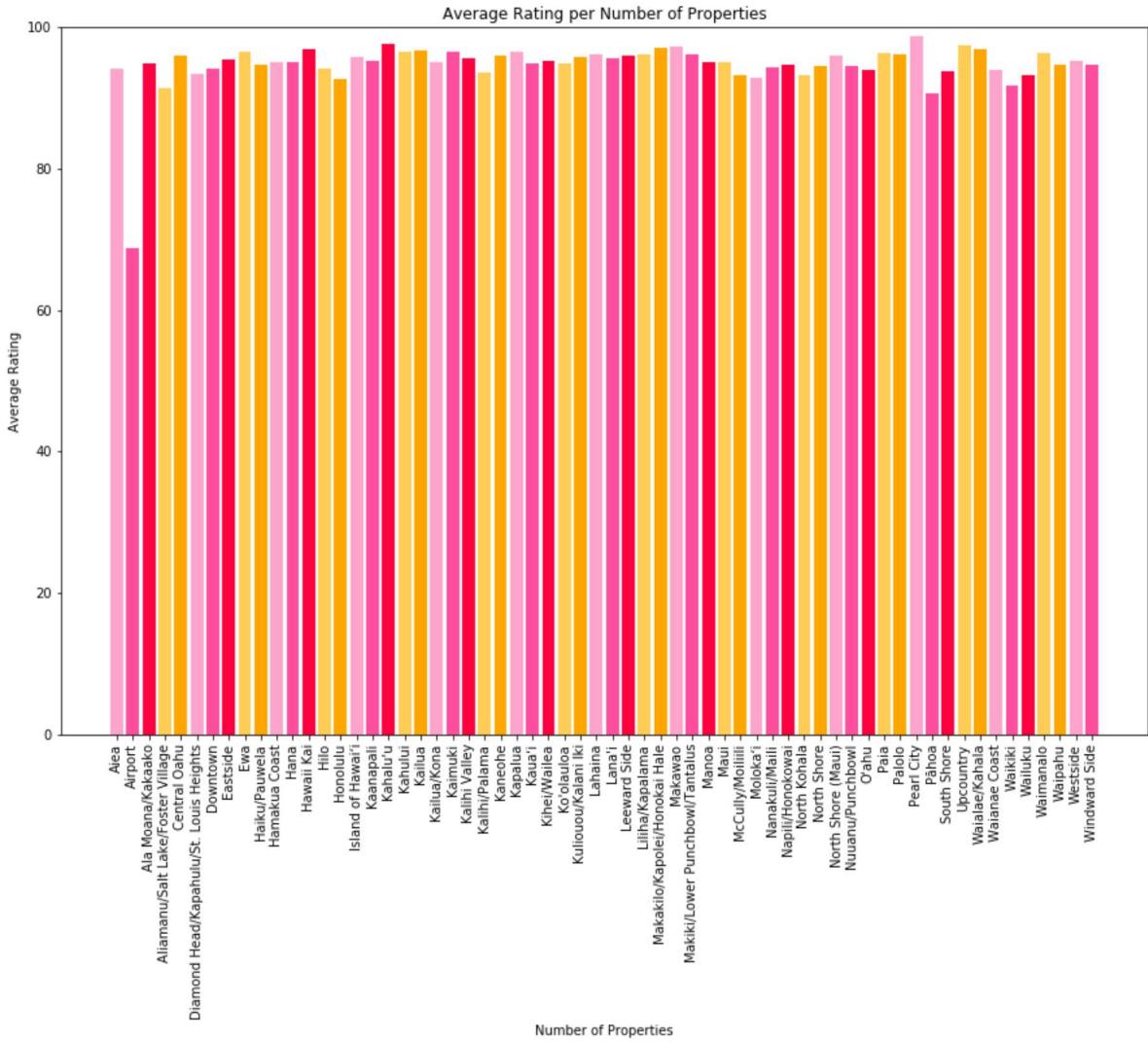
Out[20]: F_onewayResult(statistic=43.20051018106537, pvalue=1.9870577727779467e-44)

Define neighbourhood

In [21]:

```
1 neighborhood_list = df["neighbourhood"].unique()
2
3 average_rating_per_neighborhood = df.groupby("neighbourhood").mean()["review_
4 average_rating_per_neighborhood
5
6 x_axis = neighborhood_list
7 y_axis = average_rating_per_neighborhood
8
9 cleaned_neighborhood_list = [x for x in neighborhood_list if str(x) != 'nan']
10 cleaned_neighborhood_list.sort()
11 cleaned_neighborhood_list
12
13 x_axis = cleaned_neighborhood_list
14 y_axis = average_rating_per_neighborhood
15
16 hawaiian_hibiscus_hex_palette = ["#ffa4cd", "#ff4e9e", "#ff0040", "#ffcd55",
17
18 plt.subplots(figsize=(15,10))
19 plt.xticks(rotation=90)
20 plt.title("Average Rating per Number of Properties")
21 plt.xlabel("Number of Properties")
22 plt.ylabel("Average Rating")
23 plt.ylim(0, 100)
24 plt.bar(x_axis, y_axis, color = hawaiian_hibiscus_hex_palette, align='center')
25 plt.show
```

Out[21]: <function matplotlib.pyplot.show(*args, **kw)>



1 write notes here for the code below

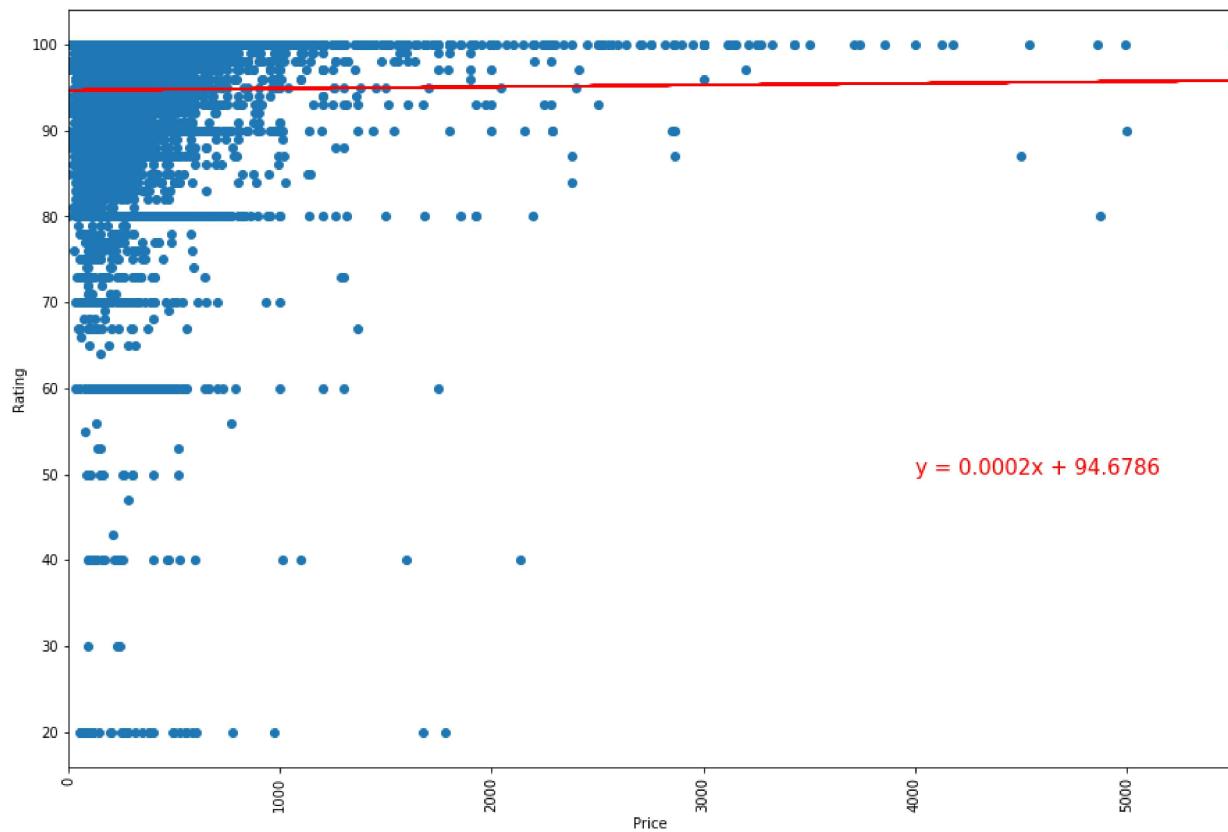
In [22]:

```

1 # Define axis
2 x_axis = df["price"]
3 y_axis = df["review_scores_rating"]
4
5
6 mask = ~np.isnan(x_axis) & ~np.isnan(y_axis)
7 slope, intercept, r_value, p_value, std_err = stats.linregress(x_axis[mask],
8 regress_values = x_axis * slope + intercept
9 line_eq = "y = " + str(round(slope,4)) + "x + " + str(round(intercept,4))
10
11 plt.subplots(figsize=(15,10))
12 plt.xticks(rotation=90)
13 plt.xlim(0, 5500)
14 plt.plot(x_axis,regress_values,"r-")
15 plt.scatter(x_axis,y_axis)
16 plt.annotate(line_eq,(4000,50),fontsize=15,color="red")
17 plt.xlabel('Price')
18 plt.ylabel('Rating')
19 print(f"The r-squared is: {r_value}")
20 plt.show()

```

The r-squared is: 0.018787686050668887



Armon Code Below: Evaluate the relationship between "number of bedrooms" and review_scores_rating.

Observations:

```
In [24]: # create a dataframe with no 0's for 'number of beds' because that would be i
          df2 = df.copy()
          df2 = df2[(df2['beds'] != 0)]
```

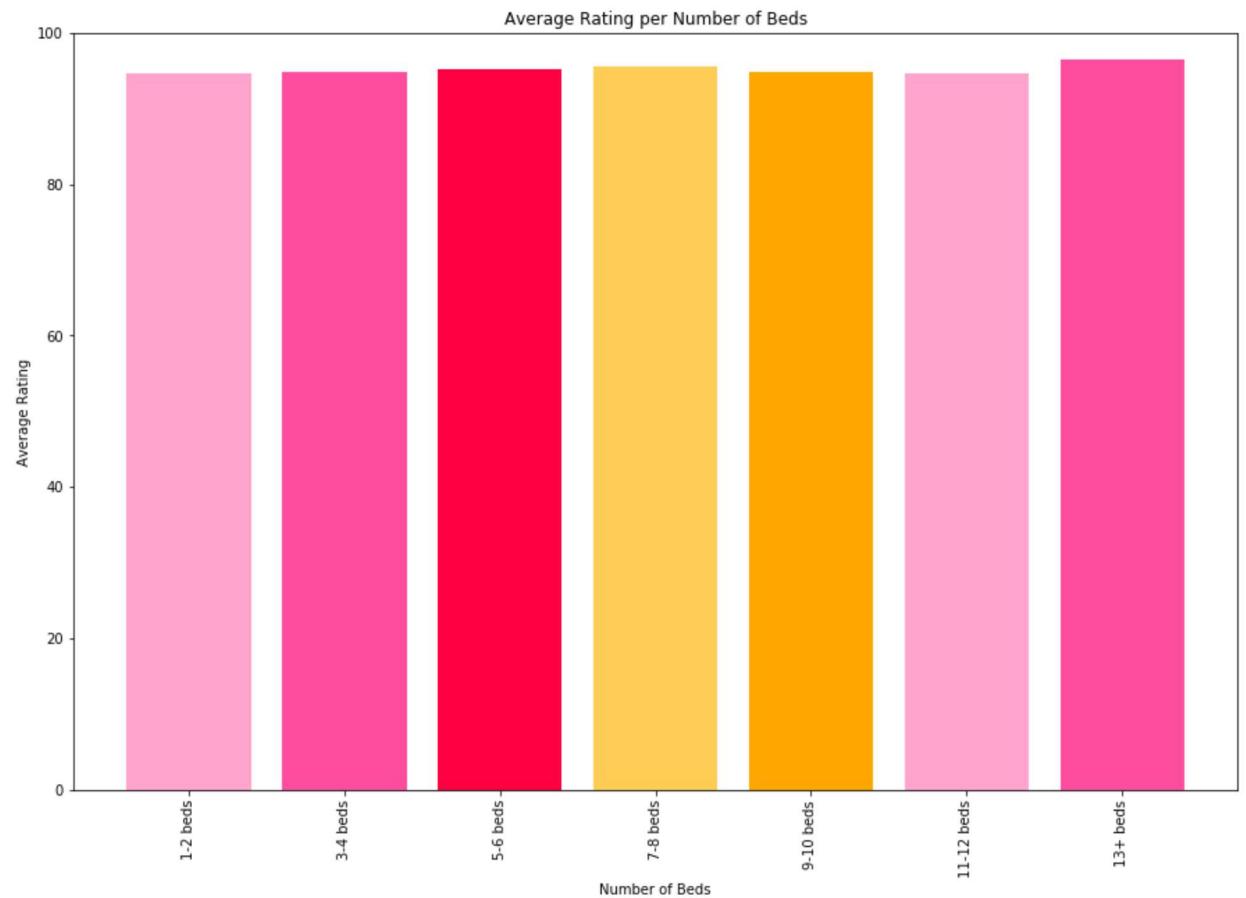


```
In [25]: # zero beds is actually an option. A 'sofa bed' is not a real bed and is not
          bins = [0, 2, 4, 6, 8, 10, 12, 20]
          group_names = ["1-2 beds", "3-4 beds", "5-6 beds", "7-8 beds", "9-10 beds", '11-12 beds", "13+ beds"]
          df2["bin_bed_count"] = pd.cut(df2["beds"], bins, labels=group_names)
          rating_mean_by_bed_count = df2.groupby("bin_bed_count").mean()["review_scores_rating"]
```

In [26]:

```
1 x_axis = group_names
2 y_axis = rating_mean_by_bed_count
3
4 hawaiian_hex_palette = ["#00739f", "#b23948", "#ff000e", "#ff407c", "#ffb2e1"
5
6 plt.subplots(figsize=(15,10))
7 plt.title("Average Rating per Number of Beds")
8 plt.xlabel("Number of Beds")
9 plt.ylabel("Average Rating")
10 plt.xticks(rotation=90)
11 plt.ylim(0, 100)
12 plt.bar(x_axis, y_axis, color = hawaiian_hibiscus_hex_palette)
13 plt.show
```

Out[26]: <function matplotlib.pyplot.show(*args, **kw)>



Lance Code Below: Evaluate the relationship between "Type of place" and review_scores_rating
Evaluate the relationship between "Island" and review_scores_rating

```
In [27]: # Define a group by every rating properties
1 average_rating_prop = df.groupby("property_type").mean()["review_scores_rating"]
#property_type = df.groupby("property_type").count()["id"]
2
3
4
5 # dropped NaN values for property types that do not have ratings in hawaii
6 average_rating_prop = average_rating_prop.dropna()
7 average_rating_prop
```

```
Out[27]: property_type
Aparthotel           92.649123
Apartment            93.820818
Barn                 95.333333
Bed and breakfast    95.806061
Boat                 95.166667
Boutique hotel       89.837209
Bungalow             95.305732
Bus                  97.000000
Cabin                94.255814
Camper/RV            93.411765
Campsite              95.500000
Casa particular (Cuba) 93.000000
Castle               96.000000
Chalet                90.200000
Condominium          94.314025
Cottage              96.919094
Dome house           95.000000
Dorm                 88.000000
Earth house           95.500000
Farm stay             96.988636
Guest suite           96.588076
Guesthouse            96.773270
Hostel                86.384615
Hotel                 90.876923
House                 95.415111
Hut                  95.000000
Island                90.714286
Loft                  94.625000
Nature lodge          88.260870
Other                 93.945205
Resort                 97.673684
Serviced apartment    94.702564
Tent                  95.631579
Timeshare              100.000000
Tiny house             95.714286
Townhouse              95.563474
Treehouse              98.307692
Vacation home          93.000000
Villa                 96.764085
Yurt                  96.200000
Name: review_scores_rating, dtype: float64
```

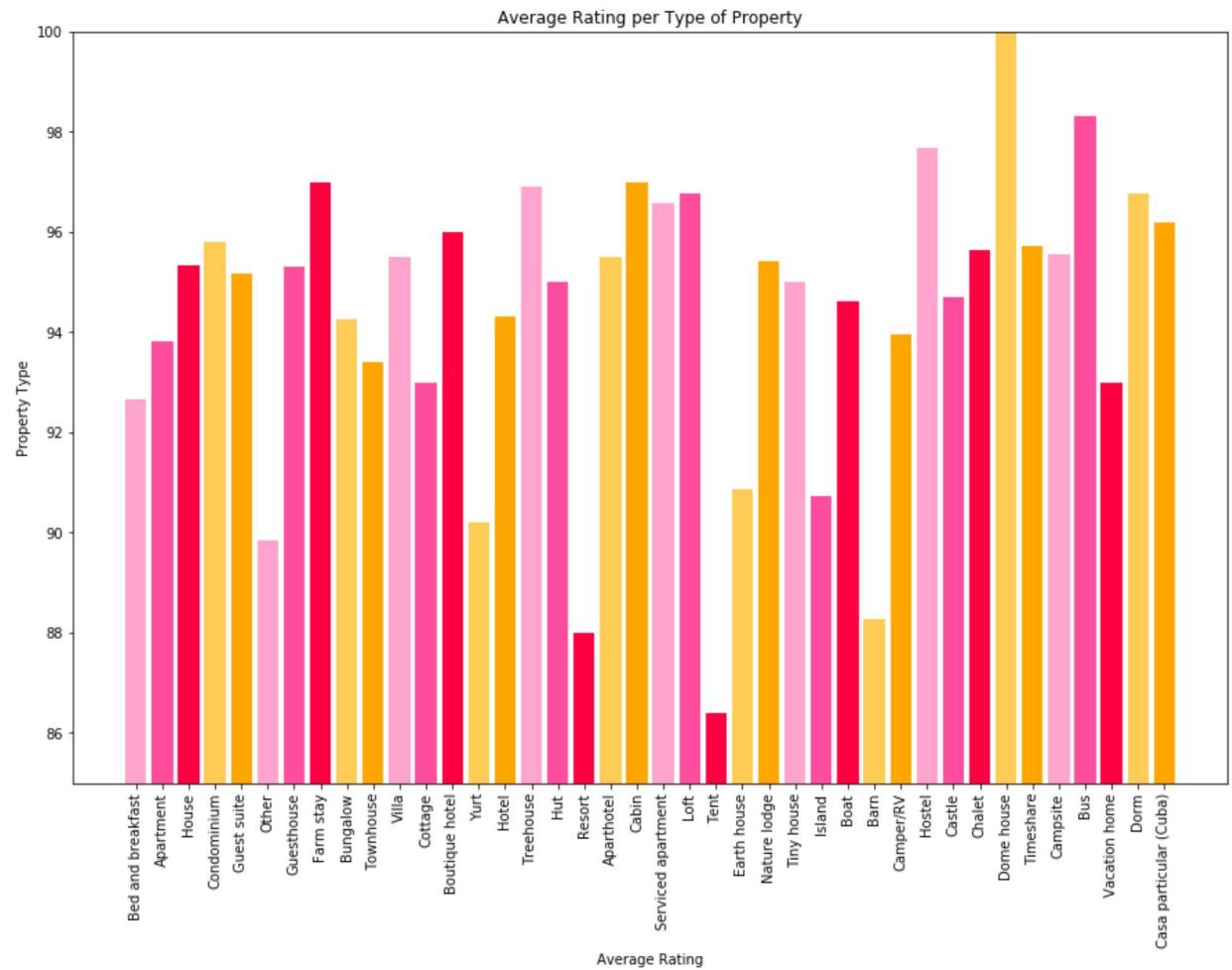
In [29]:

```

1 # set x and y axis for bar graph
2 x_axis = df["property_type"].unique()
3 y_axis = average_rating_prop
4
5 # print labels/title and x/y limits to improve visual
6 plt.subplots(figsize=(15,10))
7 plt.title("Average Rating per Type of Property")
8 plt.xlabel("Average Rating")
9 plt.ylabel("Property Type")
10 plt.xticks(rotation=90)
11 plt.ylim(85, 100)
12
13 # display graph
14 plt.bar(x_axis, y_axis, color = hawaiian_hibiscus_hex_palette)
15 plt.show

```

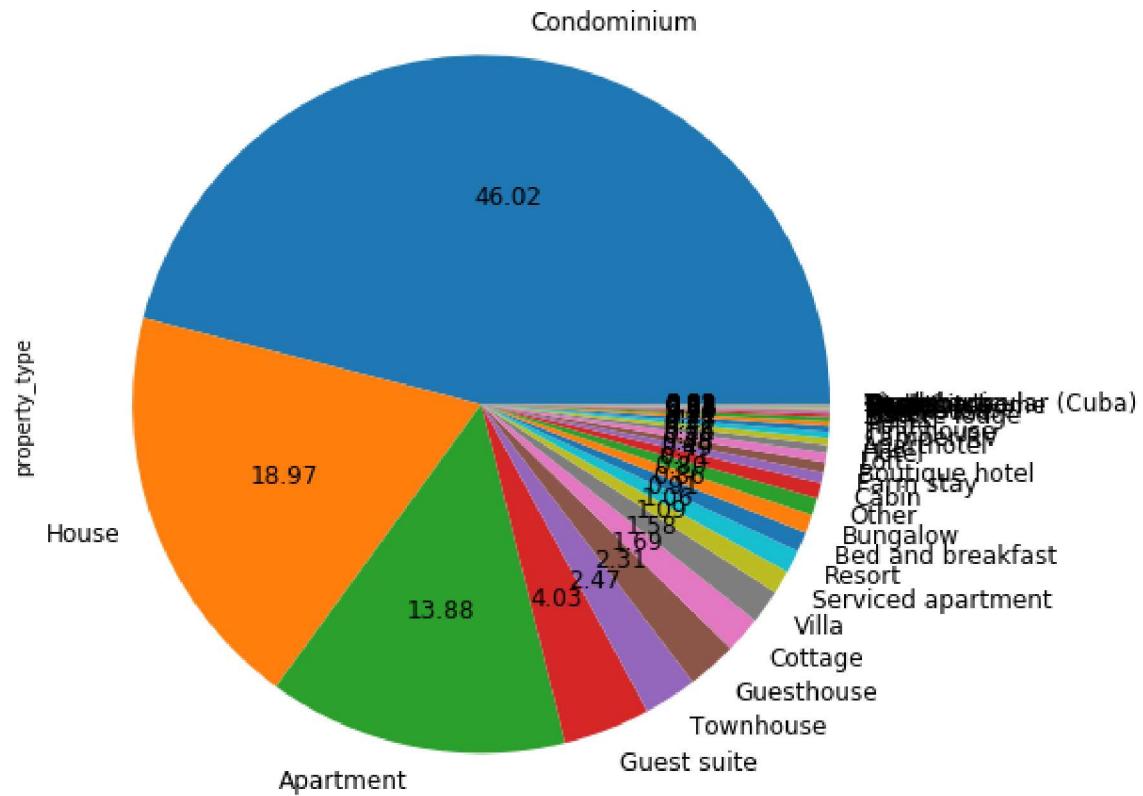
Out[29]: <function matplotlib.pyplot.show(*args, **kw)>



```
In [34]: # Build a pie to show the availability by property type in Airbnb Hawaii.  
1 prop_type_pie=df.property_type.value_counts()/len(df.property_type)  
2 prop_type_pie.plot.pie(autopct='%.2f',fontsize=12,figsize=(8,8))  
3 plt.title('Property types availability in AirBnB',fontsize=20)
```

Out[34]: Text(0.5, 1.0, 'Property types availability in AirBnB')

Property types availability in AirBnB

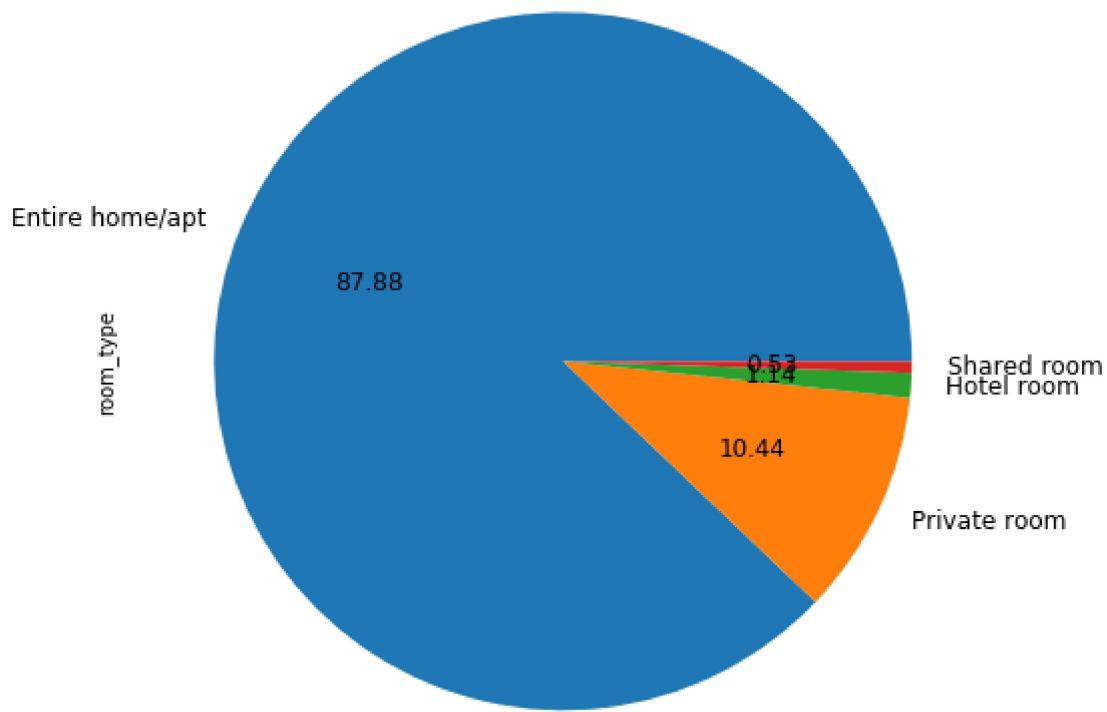


In [40]:

```
1 # Build a pie to show the availability by property room type in Airbnb Hawaii
2 room_type_pie=df.room_type.value_counts()/len(df.room_type)
3 room_type_pie.plot.pie(autopct='%.2f', fontsize=12, figsize=(8,8))
4 plt.title('Property Room types availability in AirBnB', fontsize=20)
```

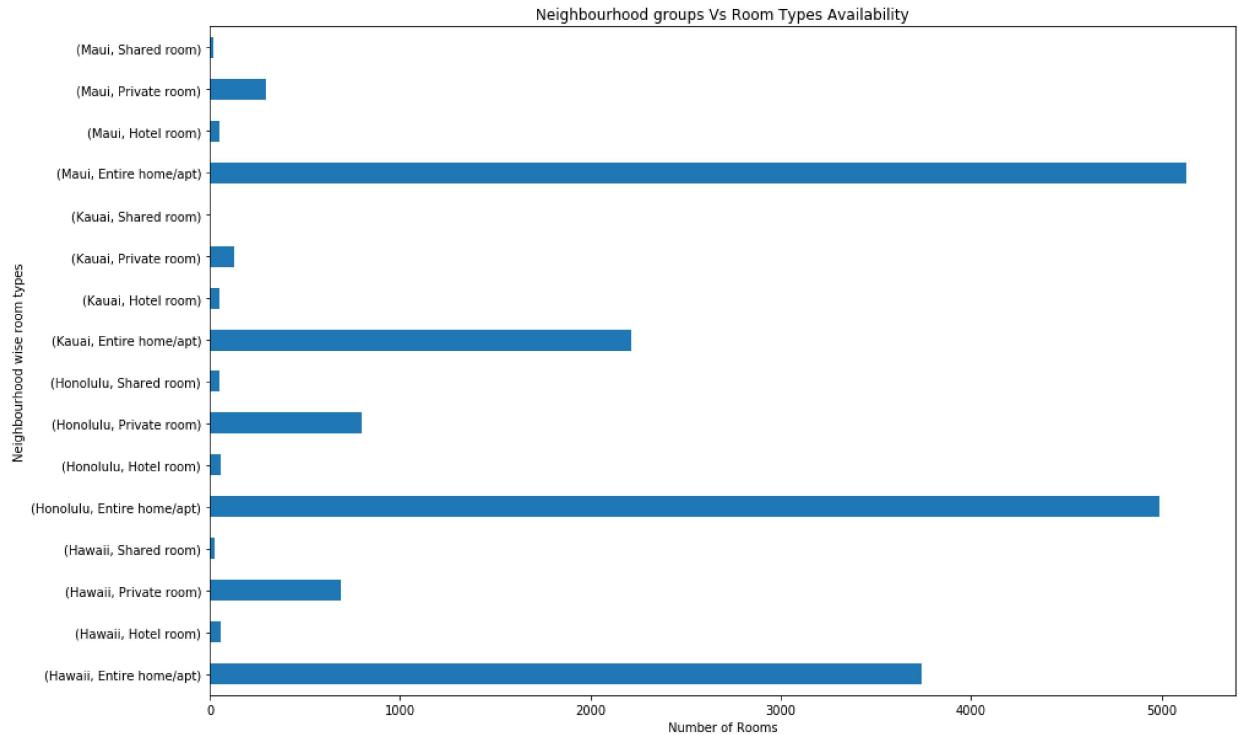
Out[40]: Text(0.5, 1.0, 'Property Room types availability in AirBnB')

Property Room types availability in AirBnB



```
In [42]: ⏎ 1 plt.subplots(figsize=(15,10))
2 df.groupby(['neighbourhood_group_cleansed', 'room_type']).room_type.count().plot()
3 plt.ylabel('Neighbourhood wise room types')
4 plt.xlabel('Number of Rooms')
5 plt.title('Neighbourhood groups Vs Room Types Availability')
```

Out[42]: Text(0.5, 1.0, 'Neighbourhood groups Vs Room Types Availability')



```
In [ ]: ⏎ 1 plt.figure(figsize=(10,6))
2 plt.scatter(df.longitude, df.latitude, c=df.availability_365, cmap='spring',
3             edgecolor='black', linewidth=1, alpha=0.75)
4
5 cbar = plt.colorbar()
6 cbar.set_label('availability_365')
```

```
1 Linear Regression
2
```

```
In [ ]: ⏎ 1
```