

# 015-assignment

May 9, 2022

Assignment: Housing in Brazil

```
[1]: import wqet_grader

wqet_grader.init("Project 1 Assessment")
```

<IPython.core.display.HTML object>

In this assignment, you'll work with a dataset of homes for sale in Brazil. Your goal is to determine if there are regional differences in the real estate market. Also, you will look at southern Brazil to see if there is a relationship between home size and price, similar to what you saw with housing in some states in Mexico.

**Note:** There are 19 graded tasks in this assignment, but you only need to complete 1.

**Before you start:** Import the libraries you'll use in this notebook: Matplotlib, pandas, and plotly. Be sure to import them under the aliases we've used in this project.

```
[2]: # Import Matplotlib, pandas, and plotly
import matplotlib.pyplot as plt
import plotly.express as px
import pandas as pd
```

## 1 Prepare Data

In this assignment, you'll work with real estate data from Brazil. In the `data` directory for this project there are two CSV that you need to import and clean.

### 1.1 Import

**Task 1.5.1:** Import the CSV file `data/brasil-real-estate-1.csv` into the DataFrame `df1`.

```
[3]: df1 = pd.read_csv("data/brasil-real-estate-1.csv")
df1.head()
```

```
[3]:  property_type  place_with_parent_names  region  lat-lon \
0      apartment  |Brasil|Alagoas|Maceió|  Northeast  -9.6443051,-35.7088142
1      apartment  |Brasil|Alagoas|Maceió|  Northeast   -9.6430934,-35.70484
2          house  |Brasil|Alagoas|Maceió|  Northeast  -9.6227033,-35.7297953
3      apartment  |Brasil|Alagoas|Maceió|  Northeast   -9.622837,-35.719556
```

```

4      apartment |Brasil|Alagoas|Maceió| Northeast      -9.654955,-35.700227

      area_m2    price_usd
0      110.0  $187,230.85
1       65.0   $81,133.37
2      211.0  $154,465.45
3       99.0  $146,013.20
4       55.0  $101,416.71

```

```
[4]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.1", df1)
```

<IPython.core.display.HTML object>

Before you move to the next task, take a moment to inspect `df1` using the `info` and `head` methods. What issues do you see in the data? What cleaning will you need to do before you can conduct your analysis?

```
[5]: df1.info()
      df1.head()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12834 entries, 0 to 12833
Data columns (total 6 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   property_type                        12834 non-null  object
1   place_with_parent_names             12834 non-null  object
2   region                              12834 non-null  object
3   lat-lon                             11551 non-null  object
4   area_m2                             12834 non-null  float64
5   price_usd                           12834 non-null  object
dtypes: float64(1), object(5)
memory usage: 601.7+ KB

```

```

[5]:  property_type  place_with_parent_names  region  lat-lon \
0      apartment |Brasil|Alagoas|Maceió| Northeast -9.6443051,-35.7088142
1      apartment |Brasil|Alagoas|Maceió| Northeast  -9.6430934,-35.70484
2        house |Brasil|Alagoas|Maceió| Northeast -9.6227033,-35.7297953
3      apartment |Brasil|Alagoas|Maceió| Northeast  -9.622837,-35.719556
4      apartment |Brasil|Alagoas|Maceió| Northeast  -9.654955,-35.700227

      area_m2    price_usd
0      110.0  $187,230.85
1       65.0   $81,133.37
2      211.0  $154,465.45
3       99.0  $146,013.20
4       55.0  $101,416.71

```

**Task 1.5.2:** Drop all rows with NaN values from the DataFrame df1.

```
[6]: #dropping rows with missing values
df1.dropna(inplace=True)
```

```
[7]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.2", df1)
```

<IPython.core.display.HTML object>

```
[8]: df1.shape
```

```
[8]: (11551, 6)
```

**Task 1.5.3:** Use the "lat-lon" column to create two separate columns in df1: "lat" and "lon". Make sure that the data type for these new columns is float.

```
[9]: #splitting lat-lon columns into lat and lon
df1[["lat", "lon"]] = df1["lat-lon"].str.split(",", expand=True)
#casting lat and lon to float
df1["lat"] = df1.lat.astype(float)
df1["lon"] = df1.lon.astype(float)

df1.head()
df1.info()
```

<class 'pandas.core.frame.DataFrame'>

Int64Index: 11551 entries, 0 to 12833

Data columns (total 8 columns):

#	Column	Non-Null Count	Dtype
0	property_type	11551 non-null	object
1	place_with_parent_names	11551 non-null	object
2	region	11551 non-null	object
3	lat-lon	11551 non-null	object
4	area_m2	11551 non-null	float64
5	price_usd	11551 non-null	object
6	lat	11551 non-null	float64
7	lon	11551 non-null	float64

dtypes: float64(3), object(5)

memory usage: 812.2+ KB

```
[10]: df1.shape
```

```
[10]: (11551, 8)
```

```
[11]: #df1 = df1.drop("lat-lon", axis="columns")
df1.head()
```

```
[11]:
```

	property_type	place_with_parent_names	region	lat-lon \
0	apartment	Brasil Alagoas Maceió	Northeast	-9.6443051,-35.7088142
1	apartment	Brasil Alagoas Maceió	Northeast	-9.6430934,-35.70484
2	house	Brasil Alagoas Maceió	Northeast	-9.6227033,-35.7297953
3	apartment	Brasil Alagoas Maceió	Northeast	-9.622837,-35.719556
4	apartment	Brasil Alagoas Maceió	Northeast	-9.654955,-35.700227

	area_m2	price_usd	lat	lon
0	110.0	\$187,230.85	-9.644305	-35.708814
1	65.0	\$81,133.37	-9.643093	-35.704840
2	211.0	\$154,465.45	-9.622703	-35.729795
3	99.0	\$146,013.20	-9.622837	-35.719556
4	55.0	\$101,416.71	-9.654955	-35.700227

```
[12]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.3", df1)
```

<IPython.core.display.HTML object>

**Task 1.5.4:** Use the "place\_with\_parent\_names" column to create a "state" column for df1. (Note that the state name always appears after "|Brasil|" in each string.)

```
[13]: #Extracting state for every house
df1["state"] = df1["place_with_parent_names"].str.split("|", expand=True)[2]
df1.head()
```

```
[13]:
```

	property_type	place_with_parent_names	region	lat-lon \
0	apartment	Brasil Alagoas Maceió	Northeast	-9.6443051,-35.7088142
1	apartment	Brasil Alagoas Maceió	Northeast	-9.6430934,-35.70484
2	house	Brasil Alagoas Maceió	Northeast	-9.6227033,-35.7297953
3	apartment	Brasil Alagoas Maceió	Northeast	-9.622837,-35.719556
4	apartment	Brasil Alagoas Maceió	Northeast	-9.654955,-35.700227

	area_m2	price_usd	lat	lon	state
0	110.0	\$187,230.85	-9.644305	-35.708814	Alagoas
1	65.0	\$81,133.37	-9.643093	-35.704840	Alagoas
2	211.0	\$154,465.45	-9.622703	-35.729795	Alagoas
3	99.0	\$146,013.20	-9.622837	-35.719556	Alagoas
4	55.0	\$101,416.71	-9.654955	-35.700227	Alagoas

```
[14]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.4", df1)
```

<IPython.core.display.HTML object>

**Task 1.5.5:** Transform the "price\_usd" column of df1 so that all values are floating-point numbers instead of strings.

```
[15]: #removing $ and "," in column price_usd
df1["price_usd"] = df1["price_usd"].str.replace("$", "")
df1["price_usd"] = df1["price_usd"].str.replace(",", "")
```

```
df1.head()
```

/tmp/ipykernel\_141/690928680.py:2: FutureWarning: The default value of regex will change from True to False in a future version. In addition, single character regular expressions will \*not\* be treated as literal strings when regex=True.

```
df1["price_usd"] = df1["price_usd"].str.replace("$", "")
```

```
[15]:
```

	property_type	place_with_parent_names	region	lat-lon	\
0	apartment	Brasil Alagoas Maceió	Northeast	-9.6443051,-35.7088142	
1	apartment	Brasil Alagoas Maceió	Northeast	-9.6430934,-35.70484	
2	house	Brasil Alagoas Maceió	Northeast	-9.6227033,-35.7297953	
3	apartment	Brasil Alagoas Maceió	Northeast	-9.622837,-35.719556	
4	apartment	Brasil Alagoas Maceió	Northeast	-9.654955,-35.700227	

	area_m2	price_usd	lat	lon	state
0	110.0	187230.85	-9.644305	-35.708814	Alagoas
1	65.0	81133.37	-9.643093	-35.704840	Alagoas
2	211.0	154465.45	-9.622703	-35.729795	Alagoas
3	99.0	146013.20	-9.622837	-35.719556	Alagoas
4	55.0	101416.71	-9.654955	-35.700227	Alagoas

```
[16]: #casting price_usd to float
df1["price_usd"] = df1.price_usd.astype(float)
```

```
[ ]: #casting price_usd to float
#df1["price_usd"] = df1.price_usd.astype(float)
#df1.info()
```

```
[17]: df1.head()
```

```
[17]:
```

	property_type	place_with_parent_names	region	lat-lon	\
0	apartment	Brasil Alagoas Maceió	Northeast	-9.6443051,-35.7088142	
1	apartment	Brasil Alagoas Maceió	Northeast	-9.6430934,-35.70484	
2	house	Brasil Alagoas Maceió	Northeast	-9.6227033,-35.7297953	
3	apartment	Brasil Alagoas Maceió	Northeast	-9.622837,-35.719556	
4	apartment	Brasil Alagoas Maceió	Northeast	-9.654955,-35.700227	

	area_m2	price_usd	lat	lon	state
0	110.0	187230.85	-9.644305	-35.708814	Alagoas
1	65.0	81133.37	-9.643093	-35.704840	Alagoas
2	211.0	154465.45	-9.622703	-35.729795	Alagoas
3	99.0	146013.20	-9.622837	-35.719556	Alagoas
4	55.0	101416.71	-9.654955	-35.700227	Alagoas

```
[18]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.5", df1)
```

<IPython.core.display.HTML object>

**Task 1.5.6:** Drop the "lat-lon" and "place\_with\_parent\_names" columns from df1.

```
[19]: df1 = df1.drop("lat-lon", axis="columns")
df1 = df1.drop("place_with_parent_names", axis="columns")
df1.head(20)
```

```
[19]:
```

	property_type	region	area_m2	price_usd	lat	lon	state
0	apartment	Northeast	110.0	187230.85	-9.644305	-35.708814	Alagoas
1	apartment	Northeast	65.0	81133.37	-9.643093	-35.704840	Alagoas
2	house	Northeast	211.0	154465.45	-9.622703	-35.729795	Alagoas
3	apartment	Northeast	99.0	146013.20	-9.622837	-35.719556	Alagoas
4	apartment	Northeast	55.0	101416.71	-9.654955	-35.700227	Alagoas
5	apartment	Northeast	56.0	75727.07	-9.614414	-35.735621	Alagoas
6	apartment	Northeast	68.0	110916.18	-9.584755	-35.662909	Alagoas
7	apartment	Northeast	187.0	249641.14	-9.658285	-35.703827	Alagoas
9	apartment	Northeast	90.0	115459.02	-9.660820	-35.702976	Alagoas
10	apartment	Northeast	137.0	361979.65	-9.663800	-35.711545	Alagoas
11	apartment	Northeast	101.0	131061.59	-9.661504	-35.702961	Alagoas
12	house	Northeast	140.0	99856.45	-9.697809	-35.893414	Alagoas
13	apartment	Northeast	250.0	466625.47	-9.629426	-35.699730	Alagoas
15	apartment	Northeast	136.0	162941.08	-9.649560	-35.737110	Alagoas
16	apartment	Northeast	145.0	200197.52	-9.660800	-35.705772	Alagoas
17	apartment	Northeast	138.0	188295.26	-9.649588	-35.708401	Alagoas
18	apartment	Northeast	175.0	405666.85	-9.650917	-35.706558	Alagoas
20	apartment	Northeast	122.0	187230.85	-9.658275	-35.705242	Alagoas
21	apartment	Northeast	98.0	118579.54	-9.660820	-35.702976	Alagoas
22	apartment	Northeast	107.0	183069.25	-9.661148	-35.700417	Alagoas

```
[ ]:
```

```
[20]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.6", df1)
```

<IPython.core.display.HTML object>

**Task 1.5.7:** Import the CSV file brasil-real-estate-2.csv into the DataFrame df2.

```
[21]: df2 = pd.read_csv("data/brasil-real-estate-2.csv")
df2.head()
```

```
[21]:
```

	property_type	state	region	lat	lon	area_m2	\
0	apartment	Pernambuco	Northeast	-8.134204	-34.906326	72.0	
1	apartment	Pernambuco	Northeast	-8.126664	-34.903924	136.0	
2	apartment	Pernambuco	Northeast	-8.125550	-34.907601	75.0	
3	apartment	Pernambuco	Northeast	-8.120249	-34.895920	187.0	
4	apartment	Pernambuco	Northeast	-8.142666	-34.906906	80.0	

```
price_brl
0 414222.98
```

```

1  848408.53
2  299438.28
3  848408.53
4  464129.36

```

```
[23]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.7", df2)
```

<IPython.core.display.HTML object>

Before you jump to the next task, take a look at `df2` using the `info` and `head` methods. What issues do you see in the data? How is it similar or different from `df1`?

```
[142]: df2.info()
df2.head()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 12833 entries, 0 to 12832
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   property_type    12833 non-null  object
 1   state            12833 non-null  object
 2   region           12833 non-null  object
 3   lat              12833 non-null  float64
 4   lon              12833 non-null  float64
 5   area_m2          11293 non-null  float64
 6   price_br1        12833 non-null  float64
dtypes: float64(4), object(3)
memory usage: 701.9+ KB

```

```
[142]:
```

	property_type	state	region	lat	lon	area_m2	\
0	apartment	Pernambuco	Northeast	-8.134204	-34.906326	72.0	
1	apartment	Pernambuco	Northeast	-8.126664	-34.903924	136.0	
2	apartment	Pernambuco	Northeast	-8.125550	-34.907601	75.0	
3	apartment	Pernambuco	Northeast	-8.120249	-34.895920	187.0	
4	apartment	Pernambuco	Northeast	-8.142666	-34.906906	80.0	

```

price_br1
0  414222.98
1  848408.53
2  299438.28
3  848408.53
4  464129.36

```

**Task 1.5.8:** Use the "price\_br1" column to create a new column named "price\_usd". (Keep in mind that, when this data was collected in 2015 and 2016, a US dollar cost 3.19 Brazilian reals.)

```
[42]:
```

```
[24]: #creating a new column
df2["price_usd"] = df2["price_brl"] /3.19
df2.head()
```

```
[24]:   property_type      state      region      lat      lon  area_m2  \
0    apartment  Pernambuco  Northeast -8.134204 -34.906326    72.0
1    apartment  Pernambuco  Northeast -8.126664 -34.903924   136.0
2    apartment  Pernambuco  Northeast -8.125550 -34.907601    75.0
3    apartment  Pernambuco  Northeast -8.120249 -34.895920   187.0
4    apartment  Pernambuco  Northeast -8.142666 -34.906906    80.0

      price_brl      price_usd
0  414222.98  129850.463950
1  848408.53  265958.786834
2  299438.28   93867.799373
3  848408.53  265958.786834
4  464129.36  145495.097179
```

```
[25]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.8", df2)
```

<IPython.core.display.HTML object>

**Task 1.5.9:** Drop the "price\_brl" column from df2, as well as any rows that have NaN values.

```
[26]: #dropping rows with missing values
df2.dropna(inplace=True)
df2.shape
```

```
[26]: (11293, 8)
```

```
[27]: #dropping price_brl

df2 = df2.drop("price_brl", axis="columns")
```

```
[28]: df2.shape
```

```
[28]: (11293, 7)
```

```
[29]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.9", df2)
```

<IPython.core.display.HTML object>

**Task 1.5.10:** Concatenate df1 and df2 to create a new DataFrame named df.

```
[30]: df= pd.concat([df1, df2])
print("df shape:", df.shape)
df.head()
```

```
df shape: (22844, 7)
```



```
[30]:
```

	property_type	region	area_m2	price_usd	lat	lon	state
0	apartment	Northeast	110.0	187230.85	-9.644305	-35.708814	Alagoas
1	apartment	Northeast	65.0	81133.37	-9.643093	-35.704840	Alagoas
2	house	Northeast	211.0	154465.45	-9.622703	-35.729795	Alagoas
3	apartment	Northeast	99.0	146013.20	-9.622837	-35.719556	Alagoas
4	apartment	Northeast	55.0	101416.71	-9.654955	-35.700227	Alagoas

```
[31]: df.head(20)
```

```
#df.shape
```

```
[31]:
```

	property_type	region	area_m2	price_usd	lat	lon	state
0	apartment	Northeast	110.0	187230.85	-9.644305	-35.708814	Alagoas
1	apartment	Northeast	65.0	81133.37	-9.643093	-35.704840	Alagoas
2	house	Northeast	211.0	154465.45	-9.622703	-35.729795	Alagoas
3	apartment	Northeast	99.0	146013.20	-9.622837	-35.719556	Alagoas
4	apartment	Northeast	55.0	101416.71	-9.654955	-35.700227	Alagoas
5	apartment	Northeast	56.0	75727.07	-9.614414	-35.735621	Alagoas
6	apartment	Northeast	68.0	110916.18	-9.584755	-35.662909	Alagoas
7	apartment	Northeast	187.0	249641.14	-9.658285	-35.703827	Alagoas
9	apartment	Northeast	90.0	115459.02	-9.660820	-35.702976	Alagoas
10	apartment	Northeast	137.0	361979.65	-9.663800	-35.711545	Alagoas
11	apartment	Northeast	101.0	131061.59	-9.661504	-35.702961	Alagoas
12	house	Northeast	140.0	99856.45	-9.697809	-35.893414	Alagoas
13	apartment	Northeast	250.0	466625.47	-9.629426	-35.699730	Alagoas
15	apartment	Northeast	136.0	162941.08	-9.649560	-35.737110	Alagoas
16	apartment	Northeast	145.0	200197.52	-9.660800	-35.705772	Alagoas
17	apartment	Northeast	138.0	188295.26	-9.649588	-35.708401	Alagoas
18	apartment	Northeast	175.0	405666.85	-9.650917	-35.706558	Alagoas
20	apartment	Northeast	122.0	187230.85	-9.658275	-35.705242	Alagoas
21	apartment	Northeast	98.0	118579.54	-9.660820	-35.702976	Alagoas
22	apartment	Northeast	107.0	183069.25	-9.661148	-35.700417	Alagoas

```
[32]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.10", df)
```

<IPython.core.display.HTML object>

<p><b>Frequent Question:</b> I can't pass this question, and I don't know what I've done wrong

<p><b>Tip:</b> In this assignment, you're working with data that's similar -

but not identical - the data used in the lessons. That means that you might need to make adjust

## 1.2 Explore

It's time to start exploring your data. In this section, you'll use your new data visualization skills to learn more about the regional differences in the Brazilian real estate market.

Complete the code below to create a `scatter_mapbox` showing the location of the properties in `df`.

```
[33]: fig = px.scatter_mapbox(
    df,
    lat="lat",
    lon="lon",
    center={"lat": -14.2, "lon": -51.9}, # Map will be centered on Brazil
    width=600,
    height=600,
    hover_data=["price_usd"], # Display price when hovering mouse over house
)

fig.update_layout(mapbox_style="open-street-map")

fig.show()
```



**Task 1.5.11:** Use the `describe` method to create a DataFrame `summary_stats` with the summary statistics for the `"area_m2"` and `"price_usd"` columns.

```
[34]: summary_stats = df[["area_m2", "price_usd"]].describe()
summary_stats
```

```
[34]:
```

	area_m2	price_usd
count	22844.000000	22844.000000
mean	115.020224	194987.315480
std	47.742932	103617.682978
min	53.000000	74892.340000
25%	76.000000	113898.770000
50%	103.000000	165697.555000
75%	142.000000	246900.880878

```
max      252.000000  525659.717868
```

```
[35]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.11", summary_stats)
```

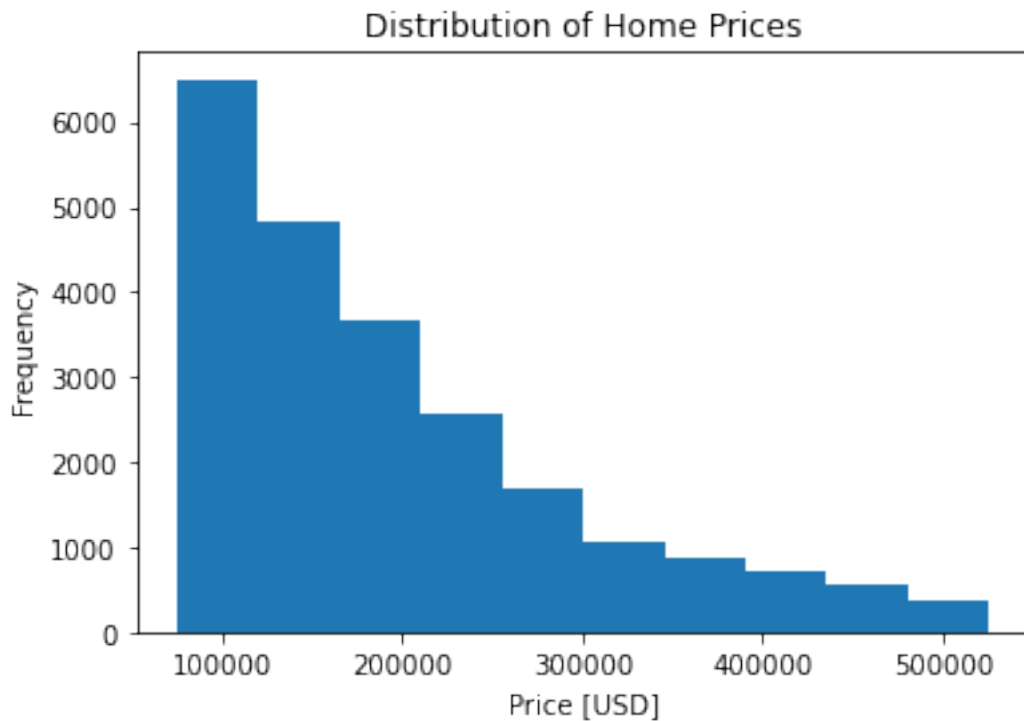
<IPython.core.display.HTML object>

**Task 1.5.12:** Create a histogram of "price\_usd". Make sure that the x-axis has the label "Price [USD]", the y-axis has the label "Frequency", and the plot has the title "Distribution of Home Prices".

```
[41]: df.head()
      df.shape
```

```
[41]: (22844, 7)
```

```
[54]: # Don't change the code below
plt.hist(df["price_usd"])
plt.xlabel("Price [USD]")
plt.ylabel("Frequency")
plt.title("Distribution of Home Prices");
plt.savefig("images/1-5-12.png", dpi=150)
```



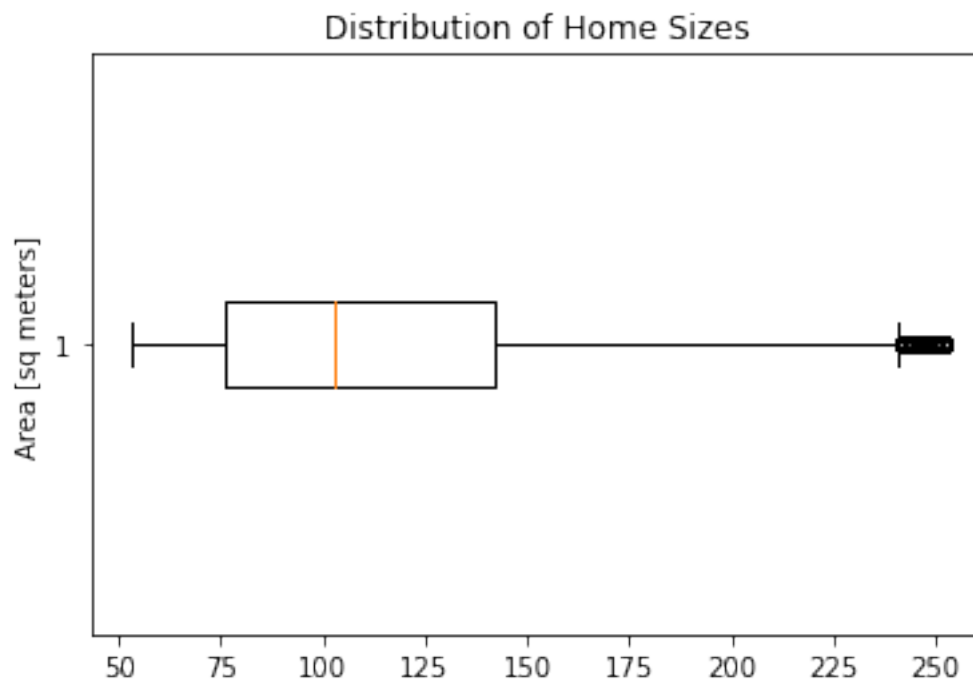
```
[ ]:
```

```
[55]: with open("images/1-5-12.png", "rb") as file:
      wqet_grader.grade("Project 1 Assessment", "Task 1.5.12", file)
```

<IPython.core.display.HTML object>

**Task 1.5.13:** Create a horizontal boxplot of "area\_m2". Make sure that the x-axis has the label "Area [sq meters]" and the plot has the title "Distribution of Home Sizes".

```
[56]: # Don't change the code below
plt.boxplot(df["area_m2"],vert=False)
plt.ylabel("Area [sq meters]")
plt.title("Distribution of Home Sizes");
plt.savefig("images/1-5-13.png", dpi=150)
```



```
[57]: with open("images/1-5-13.png", "rb") as file:
      wqet_grader.grade("Project 1 Assessment", "Task 1.5.13", file)
```

<IPython.core.display.HTML object>

**Task 1.5.14:** Use the `groupby` method to create a Series named `mean_price_by_region` that shows the mean home price in each region in Brazil, sorted from smallest to largest.

```
[58]: mean_price_by_region =df.groupby("region")["price_usd"].mean().
      ↪sort_values(ascending=False)
mean_price_by_region
```

```
[58]: region
      Southeast      208996.762778
      South          189012.345265
      Northeast      185422.985441
      North          181308.958207
      Central-West    178596.283663
      Name: price_usd, dtype: float64
```

```
[59]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.14", mean_price_by_region)
```

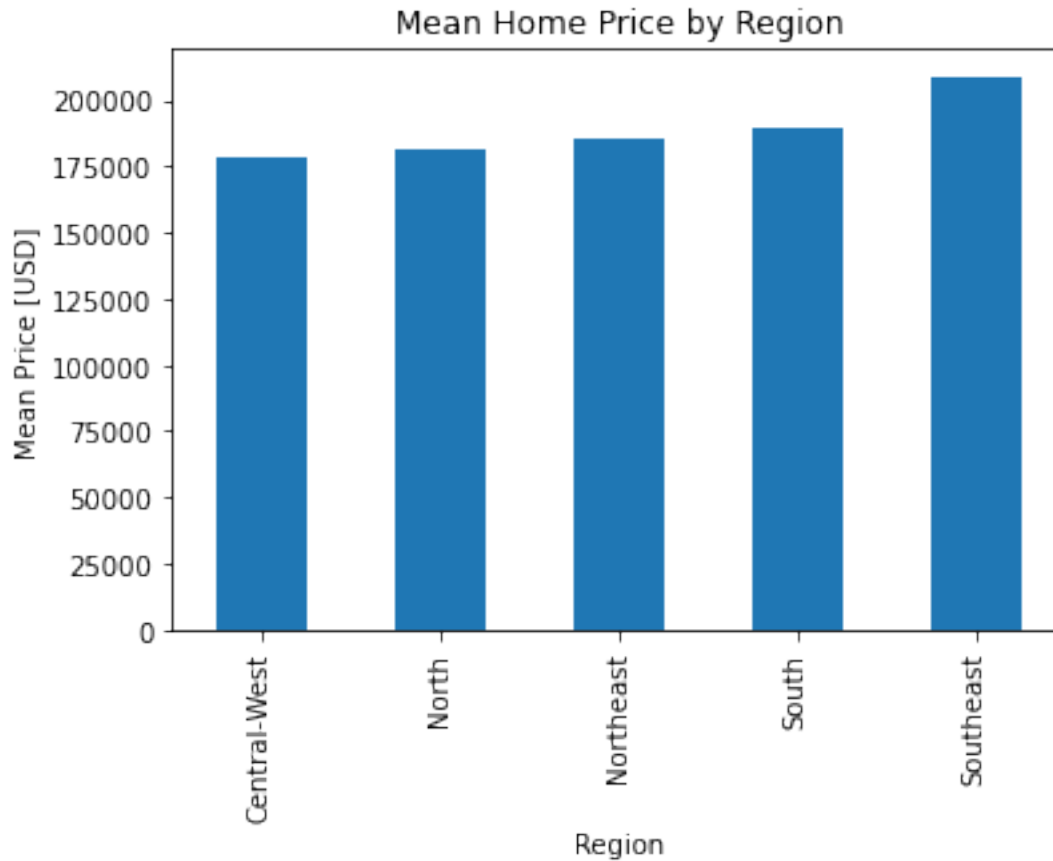
<IPython.core.display.HTML object>

```
[86]: mean_price_by_region =df.groupby("region")["price_usd"].mean()
      mean_price_by_region
```

```
[86]: region
      Central-West    178596.283663
      North          181308.958207
      Northeast      185422.985441
      South          189012.345265
      Southeast      208996.762778
      Name: price_usd, dtype: float64
```

**Task 1.5.15:** Use `mean_price_by_region` to create a bar chart. Make sure you label the x-axis as "Region" and the y-axis as "Mean Price [USD]", and give the chart the title "Mean Home Price by Region".

```
[88]: #"Mean Home Price by Region"
      # Don't change the code below
      mean_price_by_region.plot(
      kind="bar",
      xlabel="Region",
      ylabel="Mean Price [USD]",
      title="Mean Home Price by Region"
      );
      plt.savefig("images/1-5-15.png", dpi=150)
```



```
[89]: with open("images/1-5-15.png", "rb") as file:
      wqet_grader.grade("Project 1 Assessment", "Task 1.5.15", file)
```

<IPython.core.display.HTML object>

<b>Keep it up!</b> You're halfway through your data exploration. Take one last break and get ready for the final task.

You're now going to shift your focus to the southern region of Brazil, and look at the relationship between home size and price.

**Task 1.5.16:** Create a DataFrame `df_south` that contains all the homes from `df` that are in the "South" region.

```
[90]: df_south = df[df["region"]=="South"]
      df_south.head()
```

```
[90]:   property_type region  area_m2  price_usd    lat    lon  state
9304    apartment   South    127.0  296448.85 -25.455704 -49.292918  Paraná
9305    apartment   South    104.0  219996.25 -25.455704 -49.292918  Paraná
9306    apartment   South    100.0  194210.50 -25.460236 -49.293812  Paraná
9307    apartment   South     77.0  149252.94 -25.460236 -49.293812  Paraná
```

```
9308      apartment  South      73.0  144167.75 -25.460236 -49.293812  Paraná
```

```
[91]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.16", df_south)
```

<IPython.core.display.HTML object>

**Task 1.5.17:** Use the `value_counts` method to create a Series `homes_by_state` that contains the number of properties in each state in `df_south`.

```
[103]: homes_by_state = df_south["state"].value_counts().head(10)
homes_by_state.head()
```

```
[103]: Rio Grande do Sul      2643
Santa Catarina              2634
Paraná                      2544
Name: state, dtype: int64
```

```
[105]: homes_by_state.max()
```

```
[105]: 2643
```

```
[104]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.17", homes_by_state)
```

<IPython.core.display.HTML object>

**Task 1.5.18:** Create a scatter plot showing price vs. area for the state in `df_south` that has the largest number of properties. Be sure to label the x-axis "Area [sq meters]" and the y-axis "Price [USD]"; and use the title "<name of state>: Price vs. Area".

<p><b>Tip:</b> You should replace <code>&lt;name of state&gt;</code> with the name of the state</p>

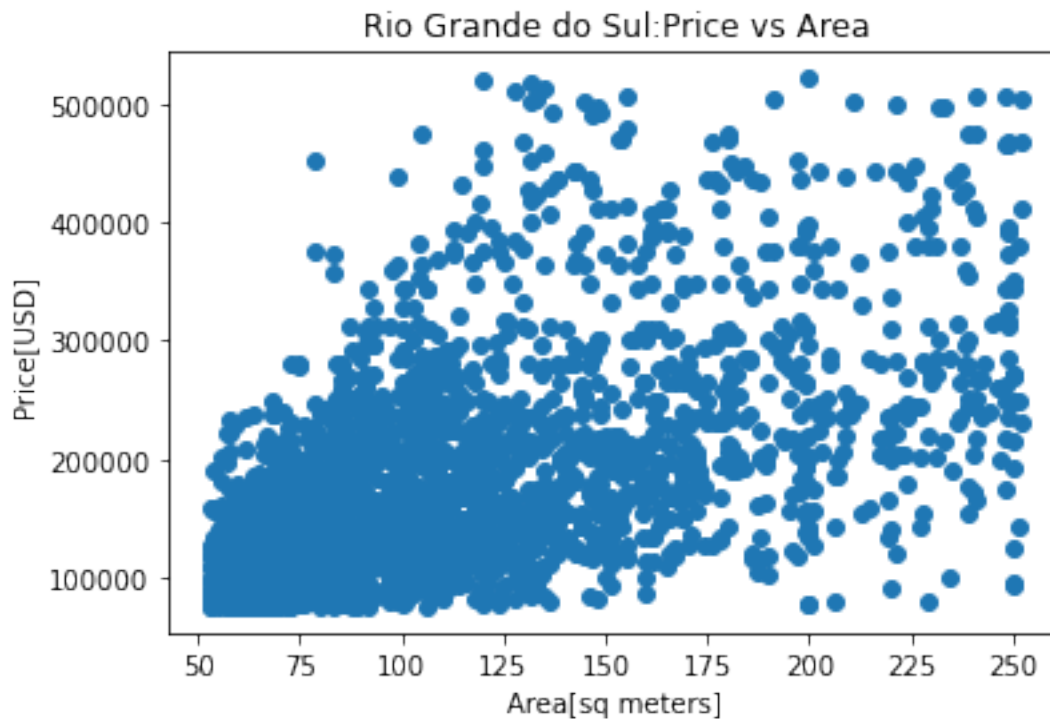
```
[109]: df_largest = df[df["state"]=="Rio Grande do Sul"]
df_largest.head()
```

```
[109]:   property_type region  area_m2    price_usd    lat    lon \
743      house  South   188.0  115770.288401 -30.027105 -51.130470
745  apartment  South    65.0  123430.141066 -30.039816 -51.223164
746  apartment  South   142.0  185145.222571 -29.696850 -53.858382
748  apartment  South   151.0  256571.996865 -30.033820 -51.198596
750  apartment  South    68.0   75957.012539 -30.034061 -51.135494

      state
743  Rio Grande do Sul
745  Rio Grande do Sul
746  Rio Grande do Sul
748  Rio Grande do Sul
750  Rio Grande do Sul
```

```
[110]: # Don't change the code below
plt.scatter(x=df_largest["area_m2"], y=df_largest["price_usd"])
plt.xlabel("Area[sq meters]")
plt.ylabel("Price[USD]")
plt.title("Rio Grande do Sul:Price vs Area");

plt.savefig("images/1-5-18.png", dpi=150)
```



```
[111]: with open("images/1-5-18.png", "rb") as file:
        wqet_grader.grade("Project 1 Assessment", "Task 1.5.18", file)
```

<IPython.core.display.HTML object>

**Task 1.5.19:** Create a dictionary `south_states_corr`, where the keys are the names of the three states in the "South" region of Brazil, and their associated values are the correlation coefficient between "area\_m2" and "price\_usd" in that state.

As an example, here's a dictionary with the states and correlation coefficients for the Southeast region. Since you're looking at a different region, the states and coefficients will be different, but the structure of the dictionary will be the same.

```
{'Espírito Santo': 0.6311332554173303,
 'Minas Gerais': 0.5830029036378931,
 'Rio de Janeiro': 0.4554077103515366,
 'São Paulo': 0.45882050624839366}
```



```
[114]: df_Rio =df[df["state"]=="Rio Grande do Sul"]
df_Santa =df[df["state"]=="Santa Catarina"]
df_Parana=df[df["state"]=="Paraná"]
```

```
[116]: rio_correlation =df_Rio["area_m2"].corr(df_Rio["price_usd"])
santa_correlation =df_Santa["area_m2"].corr(df_Santa["price_usd"])
parana_correlation =df_Parana["area_m2"].corr(df_Parana["price_usd"])
```

```
[117]: correlation_consts=[rio_correlation ,santa_correlation ,parana_correlation ]
states=["Rio Grande do Sul","Santa Catarina","Paraná"]
south_states_corr = dict(zip(states, correlation_consts))

south_states_corr
```

```
[117]: {'Rio Grande do Sul': 0.5773267433717683,
'Santa Catarina': 0.5068121776366781,
'Paraná': 0.5436659935502659}
```

```
[120]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.19", south_states_corr)
```

```
-----
Exception                                Traceback (most recent call last)
Input In [120], in <cell line: 1>()
----> 1_
↳ wqet_grader.grade("Project 1 Assessment", "Task 1.5.19", south_states_corr)

File /opt/conda/lib/python3.9/site-packages/wqet_grader/__init__.py:180, in_
↳ grade(assessment_id, question_id, submission)
    175 def grade(assessment_id, question_id, submission):
    176     submission_object = {
    177         'type': 'simple',
    178         'argument': [submission]
    179     }
--> 180     return_
↳ show_score(grade_submission(assessment_id, question_id, submission_object))

File /opt/conda/lib/python3.9/site-packages/wqet_grader/transport.py:145, in_
↳ grade_submission(assessment_id, question_id, submission_object)
    143     raise Exception('Grader raised error: {}'.format(error['message']))
    144     else:
--> 145     raise Exception('Could not grade submission: {}'.
↳ format(error['message']))
    146 result = envelope['data']['result']
    148 # Used only in testing
```

**Exception:** Could not grade submission: Could not verify access to this  
↪assessment: Received error from WQET submission API: You have already passed  
↪this course!

---

Copyright © 2022 WorldQuant University. This content is licensed solely for personal use. Redistribution or publication of this material is strictly prohibited.