

015-assignment

April 25, 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 pandas as pd
import plotly.express as px
```

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`.

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

```
[26]:
```

	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()
```

```

<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

```

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

```

[6]:  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.

```
[27]: df1.dropna(inplace = True)
```

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

<IPython.core.display.HTML object>

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.

```
[28]: df1[['lat', 'lon']] = (  
    df1['lat-lon'].str.split(',', expand = True).astype('float')  
)
```

```
[14]: 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
```

```
[15]: 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.)

```
[29]: df1['state'] = df1['place_with_parent_names'].str.split('|', expand = True)[2]
```

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

```
[20]:
```

	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

```

```
[21]: 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.

```
[30]: df1['price_usd'] = (df1['price_usd']
                          .str.replace('$', '', regex = False)
                          .str.replace(',', '')
                          .astype('float')
                          )
```

```
[31]: 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.

```
[32]: df1.drop(columns = ['lat-lon', 'place_with_parent_names'], inplace = True)
```

```
[33]: 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.

```
[34]: df2 = pd.read_csv('data/brasil-real-estate-2.csv')
df2.info()
```

```

<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

```

```
[35]: 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`?

```
[36]: df2.info()
```

```

<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

```

```
[37]: df2.head()
```

```

[37]:  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.)

```
[38]: df2['price_usd'] = df2['price_br1'] / 3.19
```

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

<IPython.core.display.HTML object>

Task 1.5.9: Drop the "price_br1" column from df2, as well as any rows that have NaN values.

```
[40]: df2.drop(columns = ['price_br1'], inplace = True)
df2.dropna(inplace = True)
```

```
[41]: 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.

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

df shape: (22844, 7)

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

<IPython.core.display.HTML object>

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

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`.

```
[44]: 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.

```
[45]: summary_stats = df[['area_m2', 'price_usd']].describe()
summary_stats
```

```
[45]:
```

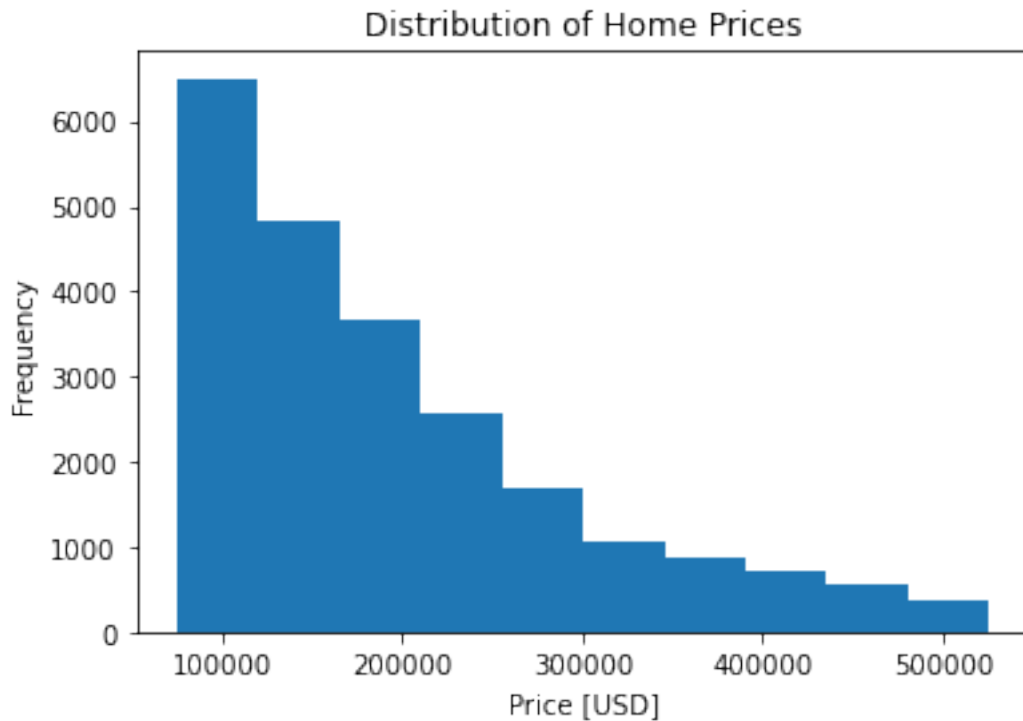
	area_m2	price_usd
count	22844.000000	22844.000000
mean	115.020224	194987.315480
std	47.742932	103617.682978
min	53.000000	74892.340000
25%	76.000000	113898.770000
50%	103.000000	165697.555000
75%	142.000000	246900.880878
max	252.000000	525659.717868

```
[46]: 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".

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

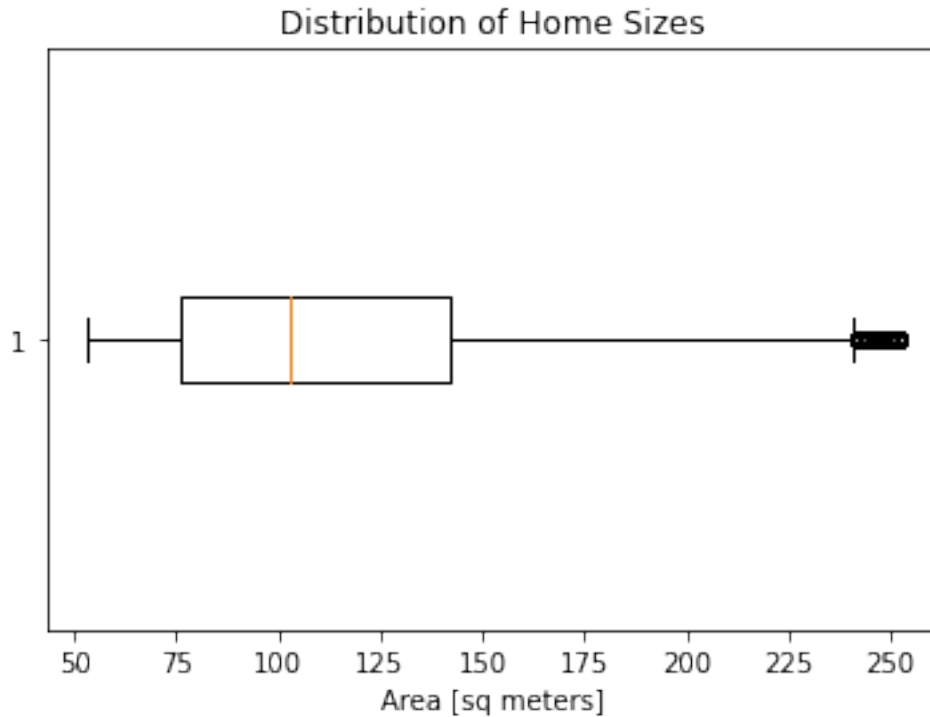


```
[50]: with open("images/1-5-12.png", "rb") as file:
      wqet_grader.grade("Project 1 Assessment", "Task 1.5.12", file)
```

<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".

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

```
[52]: 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.

```
[61]: mean_price_by_region = df.groupby('region')['price_usd'].mean().sort_values()
      mean_price_by_region
```

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

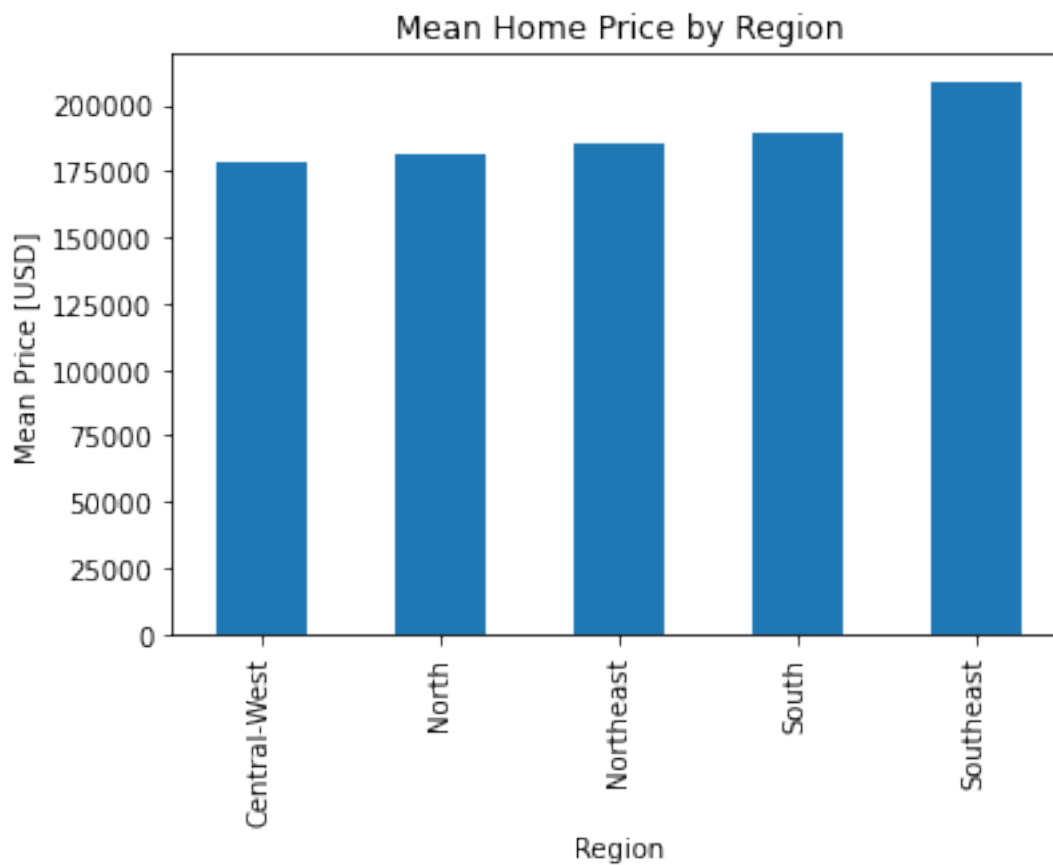
```
[56]: wqet_grader.grade("Project 1 Assessment", "Task 1.5.14", mean_price_by_region)
```

<IPython.core.display.HTML object>

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".

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



```
[63]: 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>

Keep it up! 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.

```
[64]: df_south = df[df['region'] == 'South']
df_south.head()
```

```
[64]:
```

	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

```
[65]: 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`.

```
[79]: homes_by_state = df_south[['state']].value_counts()
homes_by_state
```

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

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

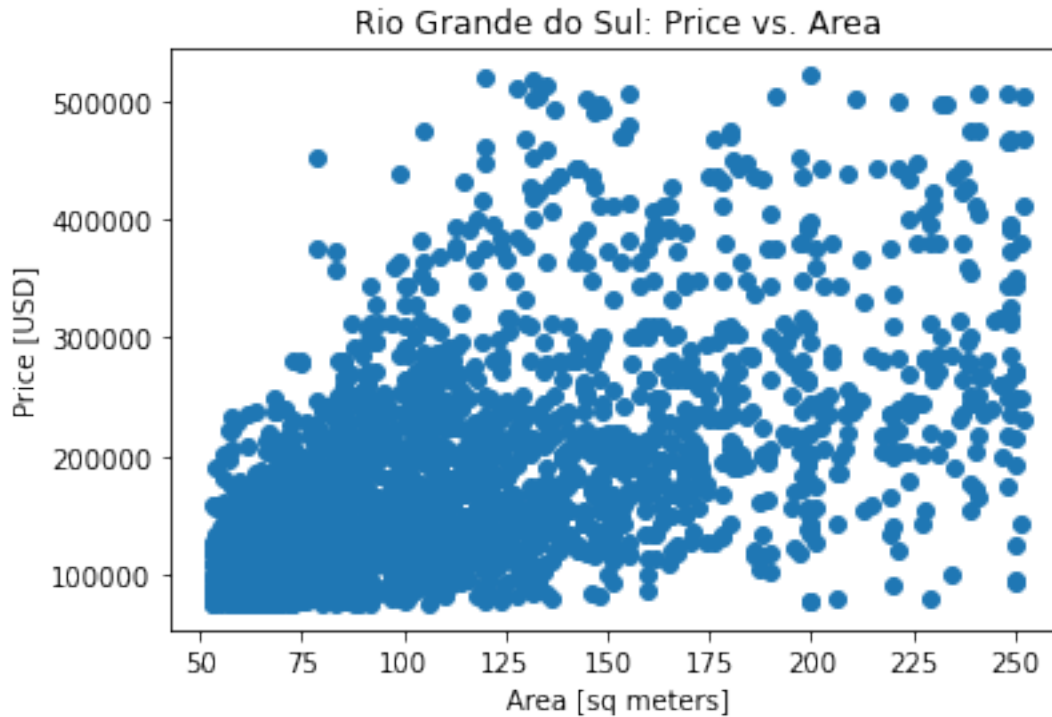
<IPython.core.display.HTML object>

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>Tip: You should replace <code><name of state></code> with the name of the state.

```
[81]: df_rio = df_south[df_south['state'] == 'Rio Grande do Sul']
plt.scatter(x = df_rio['area_m2'], y = df_rio['price_usd'])
plt.xlabel('Area [sq meters]')
plt.ylabel('Price [USD]')
plt.title('Rio Grande do Sul: Price vs. Area')

# Don't change the code below
plt.savefig("images/1-5-18.png", dpi=150)
```



```
[82]: 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}
```

```
[85]: df_south[df_south['state'] == 'Rio Grande do Sul']['area_m2'].
      ↪corr(df_south[df_south['state'] == 'Rio Grande do Sul']['price_usd'])
```

```
[85]: 0.5773267433717683
```

```
[86]: df_south[df_south['state'] == 'Santa Catarina']['area_m2'].
      ↪corr(df_south[df_south['state'] == 'Santa Catarina']['price_usd'])
```

```
[86]: 0.5068121776366781
```

```
[87]: df_south[df_south['state'] == 'Paraná']['area_m2'].  
      ↪corr(df_south[df_south['state'] == 'Paraná']['price_usd'])
```

```
[87]: 0.5436659935502659
```

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

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

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

```
-----  
Exception                                Traceback (most recent call last)  
Input In [84], 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.