# Course Title:

# Data Science (数据科学)

**(Semester: Fall 2021)**

# Dr. Oluwarotimi W. SAMUEL

**Research Center for Neural Engineering**

**Shenzhen Institutes of Advanced Technology**

**Chinese Academy of Sciences**

**Contact:** (Email: samuel@siat.ac.cn  & timitex92@gmail.com)

Phone: +86-15814491870

**(2021.10.28)**

# Data Visualization

## ❑ Outline for today's lecture

- ✓ Data Visualization

- ✓ Tools for Visualizing Data

- ✓ Case Study Using Online Dataset

# Data Visualization

❑ **Objective:** This lecture will focus on Data Visualization with practical application using a real-life data.

❑ **Expectation:** At the end of this lecture, students are expected to understand the procedure for visualizing data in the context of a data science project.

❑ **Data Visualization Definition:**

✓ Data visualization is an important step in DS process that involves the use of graphical/pictorial representation to present information/data.

✓ This allows us to better observe trends/patterns in our data set.

✓ Taking a look at a large set of numbers in a tabular form may limit understanding of inherent trends.
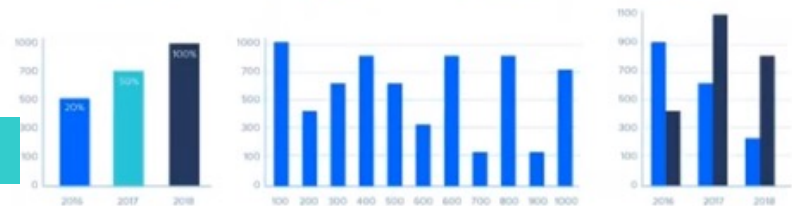
❑ **Data Description via Visualization:**

✓ Hence, visualization of data using graphs, charts, and other forms of graphical representation is essential.



**Pie chart**

**Bar plots**

**Line plots**

❑ **Importance of Data Visualization:**

| 1 | To view changes over time easily via a visual aid rather that plain data |
|---|---|
| 2 | To discover correlations among two or more variables seamlessly |
| 3 | To simplify complex information into user-friendly formats |
| 4 | To tell a better story with a bunch of pictures over time |

## ❑ **Visualization:**

✓ Note that, Data Visualization is an essential tool for Data Science.

✓ Visualization brings Data to life.

✓ Often, a good visualization can convey trends and anomalies in the data more efficiently than a written description (text/table)

# ❑ **Data Visualization:**

Visualization can be a great way to convey your predictions and conclusions to others.

❑ **Tools for Data Visualization:**

We need to use some computational tools to create Data Visualization.

Two useful visualization tools that I will introduce to you are:

✓ Python's matplotlib

✓ Seaborn libraries

## ❑ **Two Useful Data Visualization Tools:**

Matplotlib is a library in Python that allows us to create a 2D plot of our data…

Seaborn is a library in Python that based on Matplotlib that allows us to create a Multi-Dimensional plots and more advanced visualization of our data



https://matplotlib.org/index.html

https://seaborn.pydata.org/

中国科学院深圳先进技术研究院
**SHENZHEN INSTITUTE OF ADVANCED TECHNOLOGY**
**CHINESE ACADEMY OF SCIENCES**

❑ **Useful Data Visualization Tools:**

✓ Usually, different types of charts are used to visualize different types of data (depending their types)

✓ Qualitative Data (also know as categorical data)

- Subtypes are **nominal** and **ordinal** data

   **Nominal data** has no inherent order

   **Ordinal data** falls into ordered categories

# ❑ **Case Study:**

**Airbnb--**Provide data that quantifies the impact of short-term rentals on housing and residential communities; and also provides a platform to support advocacy for policies to protect our cities from the impacts of short-term rentals.

http://insideairbnb.com/about.html

# ❑ **Case Study:**



https://www.kaggle.com/dgomonov/new-york-city-airbnb-open-data

❑ **Case Study:**

**Data Description**

Since 2008, guests and hosts have used **Airbnb** to expand on traveling possibilities and present more unique, personalized way of experiencing the world. This dataset describes the listing activity and metrics in NYC, NY for 2019.

❑ **Case Study:**

**Data Description**

This data file includes all needed information to find out more about **hosts**, **geographical availability**, **prices**, necessary metrics to make **predictions** and **draw conclusions**.

*Way To Innovation*

# Case Study:

# Inspiration

- ✓ What can we learn about different hosts and areas?

- ✓ What can we learn from predictions? (ex: locations, prices, reviews, etc.)

- ✓ Which hosts are the busiest and why?

- ✓ Is there any noticeable difference of traffic among different areas and what could be the reason for it?

# Performing Data Visualization

□ **Case Study:**

First, we need to import the following needed libraries (matplolib and seaborn) to work with the Airbnb data.

```
# -*- Data Science Practicals (Analysis of the 2019 Airbnb Data for NYC) -*- #

# In this section, relevant libraries were imported
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sn
```

# Performing Data Visualization

❑ **Case Study:**

Next, we read the csv file that contains the data using the "read_csv()" function provided by the Pandas library.

```
# Read the CSV file that contains the dataset using read_csv function
listings = pd.read_csv('Airbnb_NYC_2019.csv', low_memory=False)
```

# Performing Data Visualization

❑ **Case Study:**

Below is what the Airbnb data for NYC looks like using the following command for the display.

```
13    # Display the Airbnb dataset
14    listings
```

| | id | name | host_id | host_name | neighbourhood_group | neighbourhood | latitude | longitude | room_type | price | minimum_nights |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2539 | Clean & quiet apt home by the park | 2787 | John | Brooklyn | Kensington | 40.64749 | -73.97237 | Private room | 149 | 1 |
| 1 | 2595 | Skylit Midtown Castle | 2845 | Jennifer | Manhattan | Midtown | 40.75362 | -73.98377 | Entire home/apt | 225 | 1 |
| 2 | 3647 | THE VILLAGE OF HARLEM....NEW YORK ! | 4632 | Elisabeth | Manhattan | Harlem | 40.80902 | -73.94190 | Private room | 150 | 3 |
| 3 | 3831 | Cozy Entire Floor of Brownstone | 4869 | LisaRoxanne | Brooklyn | Clinton Hill | 40.68514 | -73.95976 | Entire home/apt | 89 | 1 |
| 4 | 5022 | Entire Apt: Spacious Studio/Loft by | 7192 | Laura | Manhattan | East Harlem | 40.79851 | -73.94399 | Entire home/apt | 80 | 10 |

# Performing Data Visualization

□ **Case Study:**

Let's say we want to create a **bar chart** that displays the count of Airbnb listings in each neighborhood group of NYC from the listings DataFrame

To do this, we will use the "countplot()" method from the "**seaborn**" library followed by the "**show**" method from "***matplolib.pyplot***"

# ❑ **Case Study:**

The result is shown as follows (count of neighborhood group):

```
17  #Let's say we want to create a bar chart that displays the count of Airbnb
18  # listings in each neighborhood group of NY from the listings DataFrame
19  # To do this, we will use the "countplot()" function from the "seaborn" library
20  # followed by the "show" method from "matplolib.pyplot"
21  sn.set_theme(style="darkgrid")
22  sn.countplot(x='neighbourhood_group', data = listings)
23  plt.show()
24
```

# Performing Data Visualization

❑ **Case Study:**

Now, let's say we want to display the **average number of** in each neighborhood group of NYC from the listings DataFrame

To do this, we will use the "**barplot()**" function from the "**seaborn**" library followed by the "**show**" method from "matplolib.pyplot"

中国科学院深圳先进技术研究院
SHENZHEN INSTITUTE OF ADVANCED TECHNOLOGY
CHINESE ACADEMY OF SCIENCES

# **Performing Data Visualization**

## ❑ **Case Study:**

Below is the results:

```
#Now, let's say we want to display the average price Airbnb listings
# in each neighborhood group of NYC from the listings DataFrame
# To do this, we will use the "barplot()" function from the "seaborn" library
# followed by the "show" method from "matplolib.pyplot"
sn.barplot(x = 'neighbourhood_group', y = 'price', data = listings)  listings: {DataFrame: (48895, 16)}
plt.show()
```

# Performing Data Visualization

☐ **Case Study:**

The **black lines** (**confidence intervals**) at the middle of each bar denotes the confidence interval that were generated by the bar plot method.

Though **confidence intervals** are an important concept, but it is not within the scope of this lecture, so we will proceed to remove it from the plot as follows.

We can achieve this by setting the "ci" parameter to False.

# Performing Data Visualization

## ❑ Case Study:

Let's say we want to reproduce the barplot above without the black lines (**confidence intervals: ci**).

Though confidence intervals are an important concept, but it is not within the scope of this lecture, so we will proceed to remove it from the plot as follows.
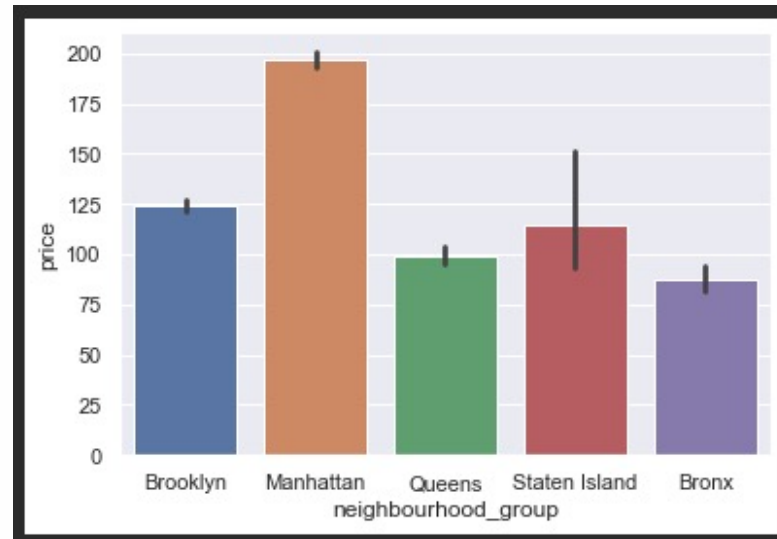
# Performing Data Visualization

## ❑ Case Study:

Below is the results:

```
#Now, let's say we want to display the average price Airbnb listings
# in each neighborhood group of NY from the listings DataFrame
# without the black lines (confidence intervals) at the middle of each bar
# We need to set the parameter "ci" to False as follows:
sn.barplot(x = 'neighbourhood_group', y = 'price', data = listings, ci = False)  listings: {DataFrame: (48895, 16)}
plt.show()
```
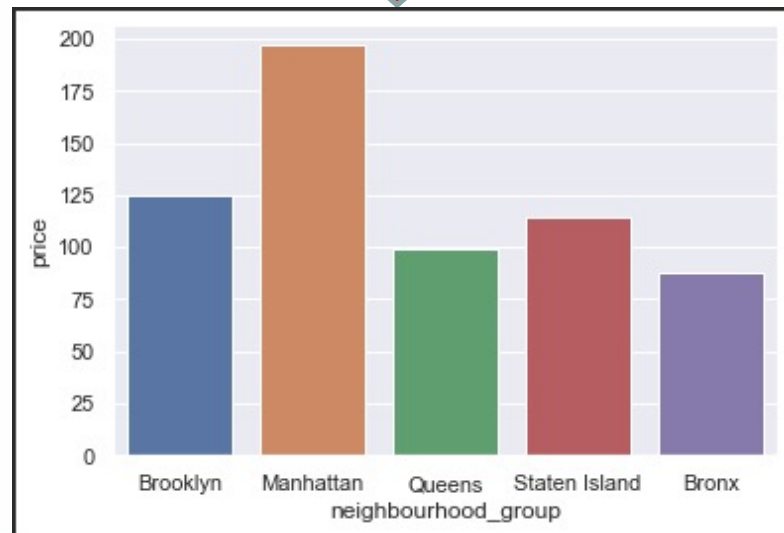
# Performing Data Visualization

❑ **Case Study:**

Below is the results:

```
#Now, let's say we want to display the average price Airbnb listings
# in each neighborhood group of NY from the listings DataFrame
# without the black lines (confidence intervals) at the middle of each bar
# We need to set the parameter "ci" to False as follows:
sn.barplot(x = 'neighbourhood_group', y = 'price', data = listings, ci = False)   listings: {DataFrame: (48895, 16)}
plt.show()
```
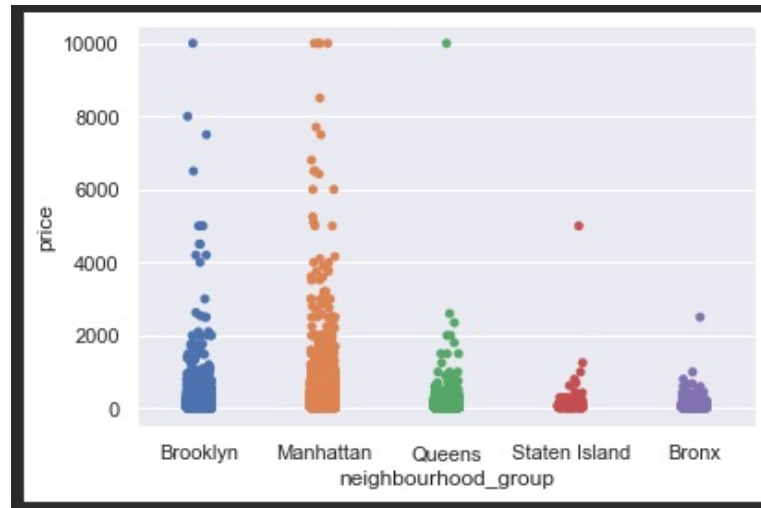
# Performing Data Visualization

## ❑ Case Study:

Below is the results using "stripplot()" function

```
#Now, let's say we want to display the average price Airbnb listings
# in each neighborhood group of NY from the listings DataFrame
# without the black lines (confidence intervals) at the middle of each bar
# We need to set the parameter "ci" to False as follows:
# sn.barplot(x = 'neighbourhood_group', y = 'price', data = listings, ci =
sn.stripplot(x = 'neighbourhood_group', y = 'price', data = listings)  list
# sn.violinplot(x = 'neighbourhood_group', y = 'price', data = listings)
```
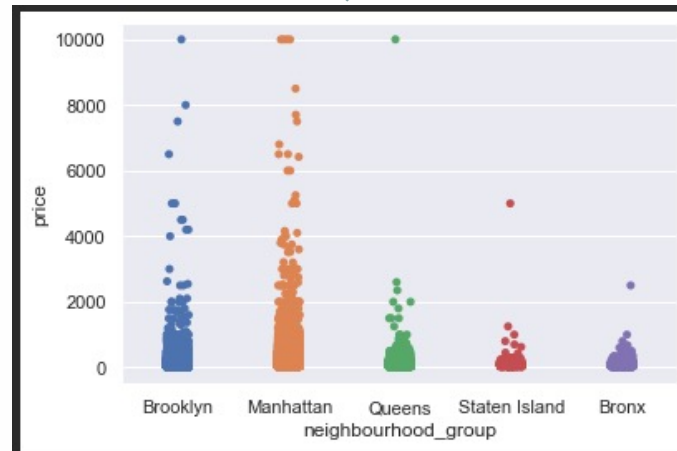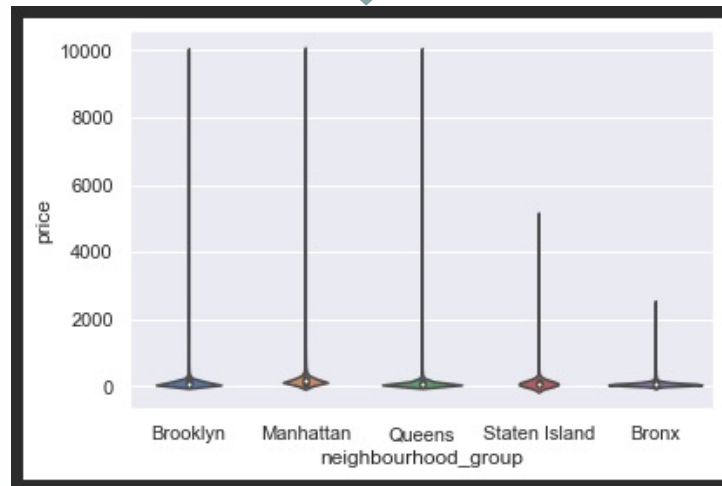
# Performing Data Visualization

## ❑ Case Study:

Below is the results using "violinplot()" function

```
#Now, let's say we want to display the average price Airbnb listings
# in each neighborhood group of NY from the listings DataFrame
# without the black lines (confidence intervals) at the middle of each bar
# We need to set the parameter "ci" to False as follows:
# sn.barplot(x = 'neighbourhood_group', y = 'price', data = listings, ci = False)
# sn.stripplot(x = 'neighbourhood_group', y = 'price', data = listings)
sn.violinplot(x = 'neighbourhood_group', y = 'price', data = listings)   listings: {DataFrame: (48895, 16)}
```

# Performing Data Visualization

☐ **Case Study:**

Next, we shall look at some common charts for visualizing Quantitative Data such as the **Histogram** and **Scatter Plots.**



**NB:** Quantitative data is another term for numerical data

# Performing Data Visualization

## ❑ Case Study:

Next, we will use the "plt.hist()" function to create a histogram of the quantitative data in the price column of the listings DataFrame
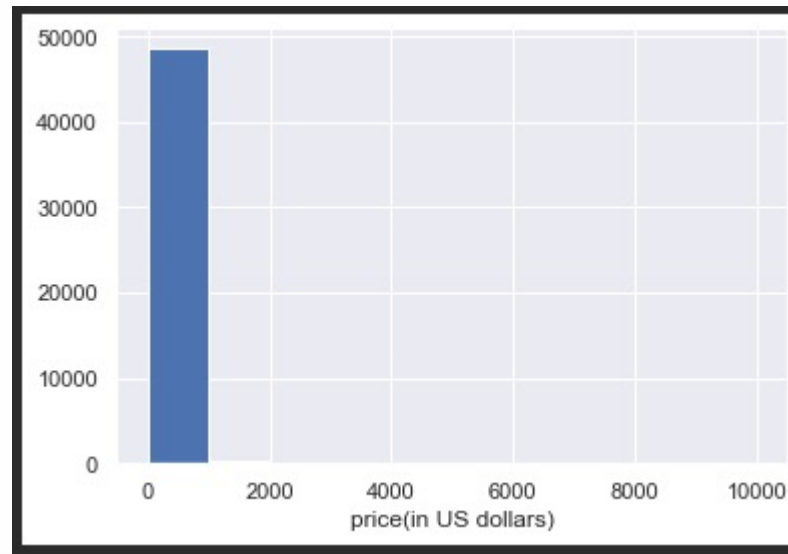
# Performing Data Visualization

❑ **Case Study:**

Below is the results:

```
# We shall use the "plt.hist()" function to create a histogram of the quantitative
# data in the price column of the listings Dataframe
plt.hist(listings['price'])
plt.xlabel('price(in US dollars)')
```
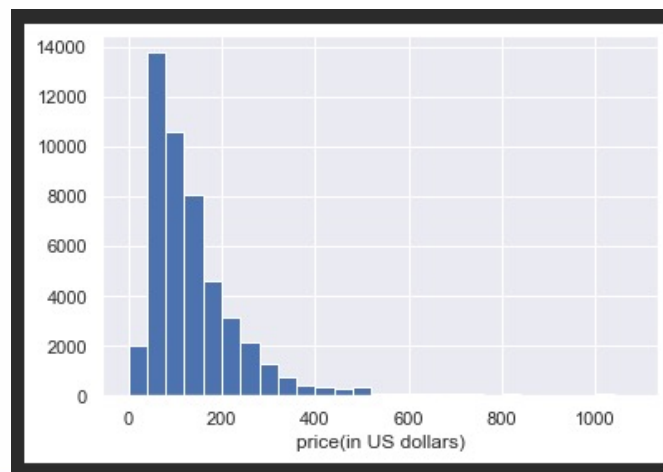
# Performing Data Visualization

❑ **Case Study:**

To better see where the data lies in the histogram, I will introduce the "bins" parameters:

```
# For better visualization of the histogram, I will introduce the 'bins' parameter as follows
plt.hist(listings['price'], bins = np.arange(0,1100,40))   listings: {DataFrame: (48895, 16)}
plt.xlabel('price(in US dollars)')
plt.show()
```

Now it is easier to see where the data lies with respect to the bins.



**A lot of the Airbnb price listings lies between 50 and 500 USD**

# **Performing Data Visualization**

❑ **Case Study:**

Next, let's talk about the **Scatter Plot**. They are often used to compare two set of quantitative data.

Assuming we want to compare the **prices** of Airbnb listings and the numbers of **reviews** for the listings.
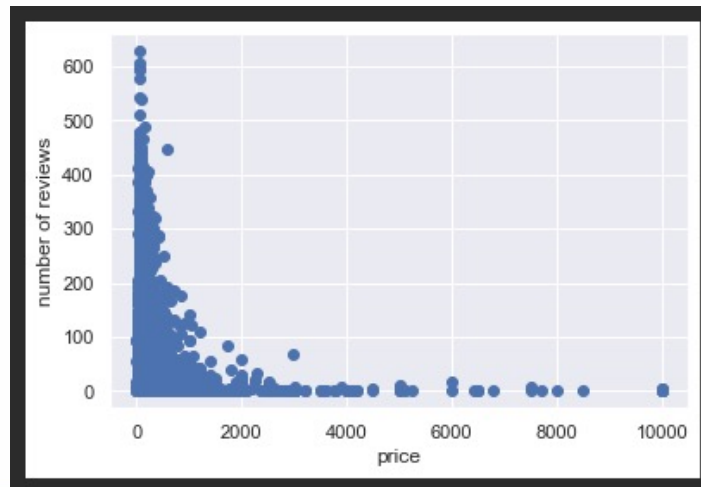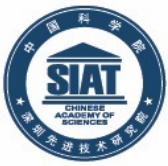
We can achieve this as follows:

# Performing Data Visualization

❑ **Case Study:**

To do this, we will use the '**scatter()**' function from the matplotlib.pyplot.

```
# For better visualization of the histogram, I will introduce the 'bins' parameter as follows
plt.scatter(x=listings ['price'], y= listings ['number_of_reviews'])   listings: {DataFrame: (48895, 16)}
plt.xlabel('price')
plt.ylabel('number of reviews')
plt.show()
```
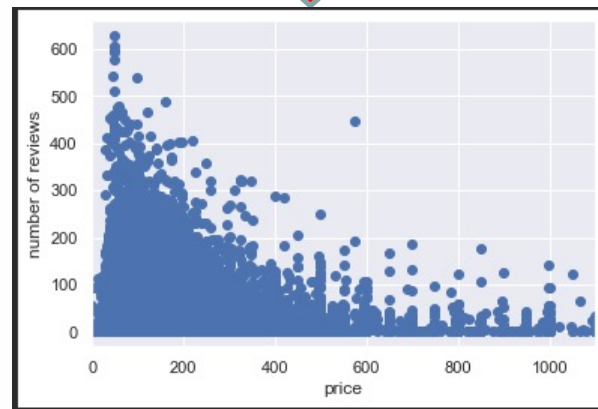
# Performing Data Visualization

## ❑ Case Study:

Let's say we want to restrict the **x-axis**, such that the scatter plot only goes up to a price of **1100**. To do this, we shall recreate the plot as follows:

```python
# Let's say I want to restrict the x-axis, such that the scatter plot only goes up to a price of 1100.
# To do this, we shall recreate the scatter plot as follows
plt.scatter(x=listings ['price'], y= listings ['number_of_reviews'])   listings: {DataFrame: (48895, 16)}
plt.xlabel('price')
plt.ylabel('number of reviews')
plt.xlim(0,1100)
plt.show()
```
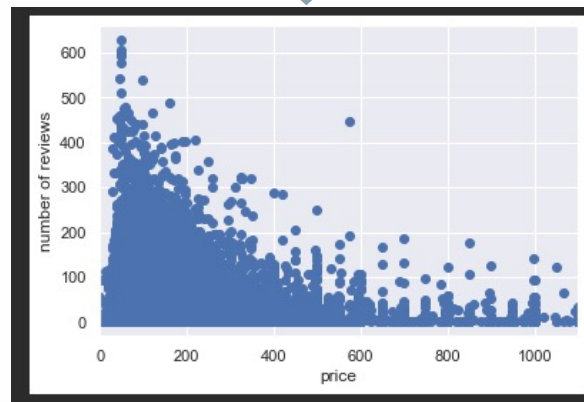
# Performing Data Visualization

## ❑ Case Study:

Let's say I want to decrease the size of points on the scatter plot To do this, we shall recreate the scatter plot using the 'scatter' method with the parameter 's' set to 5.

```
# Let's say I want to decrease the size of points on the scatter plot.
#To do this, we shall recreate the scatter plot using the 'scatter' method with the parameter 's' set to 50.
plt.scatter(x=listings ['price'], y= listings ['number_of_reviews'], s = 5)  listings: {DataFrame: (48895, 16)}
plt.xlabel('price')
plt.ylabel('number of reviews')
plt.xlim(0,1100)
plt.show()
```

# Performing Data Visualization

## ❑ Case Study:

```
# Let's say I want to decrease the size of points on the scatter plot.
#To do this, we shall recreate the scatter plot using the 'scatter' method with the parameter 's' set to 50.
plt.scatter(x=listings ['price'], y= listings ['number_of_reviews'], s = 5)  listings: {DataFrame: (48895, 16)}
plt.xlabel('price')
plt.ylabel('number of reviews')
plt.xlim(0,1100)
plt.show()
```
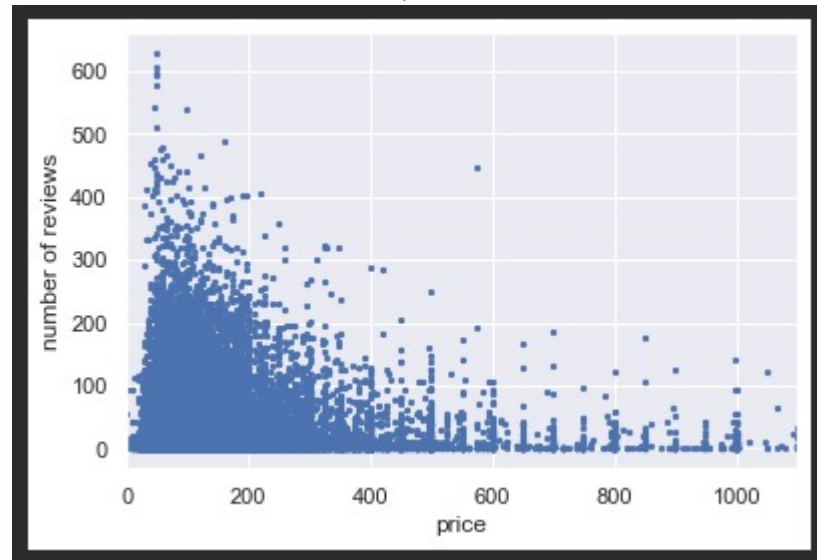
The points on this scatter plot are a bit smaller, it is a little bit easier to look at.

❑ **Performing Data Visualization:**

✓ Looking at the plot at first instance, one could observe some obvious trend…

✓ That is, listings with Lower Prices have more reviews.

✓ This leads me to ask the following questions?

✓ Is there an association between Price and the Number of Reviews for Listings?

✓ This is just one type of the questions that could come up when you are Visualizing your Data

✓ Thus, this may prompted to conduct further analysis to investigate this question.

# ❑ **Performing Data Visualization:**

✓ In summary, Visualizing your data helps you observe trends which leads you to ask questions, that then inform the kind of further analysis you will perform on the data.

✓ Now you know how to know use Histogram and Scatter plots to visualize your data.

# Questions and Comments!