

Introduction to Seaborn

About the Data

- In this notebook, we will be working with 2 datasets:
- Facebook's stock price throughout 2018 (obtained using the `stock_analysis` package)
 - Earthquake data from September 18, 2018 - October 13, 2018 (obtained from the US Geological Survey (USGS) using the USGS API)

Setup

```
1 %matplotlib inline
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import pandas as pd
5 import requests as req
6 import time as pt
7
8 # ft = plt.rcParams
9 # %matplotlib inline
10 # %matplotlib widget
11 # %matplotlib qt
12 # %matplotlib qtconsole
13 # %matplotlib qtconsole
14 # %matplotlib qtconsole
15 # %matplotlib qtconsole
16 # %matplotlib qtconsole
17 # %matplotlib qtconsole
18 # %matplotlib qtconsole
19 # %matplotlib qtconsole
20 # %matplotlib qtconsole
21 # %matplotlib qtconsole
22 # %matplotlib qtconsole
23 # %matplotlib qtconsole
24 # %matplotlib qtconsole
25 # %matplotlib qtconsole
26 # %matplotlib qtconsole
27 # %matplotlib qtconsole
28 # %matplotlib qtconsole
29 # %matplotlib qtconsole
30 # %matplotlib qtconsole
31 # %matplotlib qtconsole
32 # %matplotlib qtconsole
33 # %matplotlib qtconsole
34 # %matplotlib qtconsole
35 # %matplotlib qtconsole
36 # %matplotlib qtconsole
37 # %matplotlib qtconsole
38 # %matplotlib qtconsole
39 # %matplotlib qtconsole
40 # %matplotlib qtconsole
41 # %matplotlib qtconsole
42 # %matplotlib qtconsole
43 # %matplotlib qtconsole
44 # %matplotlib qtconsole
45 # %matplotlib qtconsole
46 # %matplotlib qtconsole
47 # %matplotlib qtconsole
48 # %matplotlib qtconsole
49 # %matplotlib qtconsole
50 # %matplotlib qtconsole
51 # %matplotlib qtconsole
52 # %matplotlib qtconsole
53 # %matplotlib qtconsole
54 # %matplotlib qtconsole
55 # %matplotlib qtconsole
56 # %matplotlib qtconsole
57 # %matplotlib qtconsole
58 # %matplotlib qtconsole
59 # %matplotlib qtconsole
60 # %matplotlib qtconsole
61 # %matplotlib qtconsole
62 # %matplotlib qtconsole
63 # %matplotlib qtconsole
64 # %matplotlib qtconsole
65 # %matplotlib qtconsole
66 # %matplotlib qtconsole
67 # %matplotlib qtconsole
68 # %matplotlib qtconsole
69 # %matplotlib qtconsole
70 # %matplotlib qtconsole
71 # %matplotlib qtconsole
72 # %matplotlib qtconsole
73 # %matplotlib qtconsole
74 # %matplotlib qtconsole
75 # %matplotlib qtconsole
76 # %matplotlib qtconsole
77 # %matplotlib qtconsole
78 # %matplotlib qtconsole
79 # %matplotlib qtconsole
80 # %matplotlib qtconsole
81 # %matplotlib qtconsole
82 # %matplotlib qtconsole
83 # %matplotlib qtconsole
84 # %matplotlib qtconsole
85 # %matplotlib qtconsole
86 # %matplotlib qtconsole
87 # %matplotlib qtconsole
88 # %matplotlib qtconsole
89 # %matplotlib qtconsole
90 # %matplotlib qtconsole
91 # %matplotlib qtconsole
92 # %matplotlib qtconsole
93 # %matplotlib qtconsole
94 # %matplotlib qtconsole
95 # %matplotlib qtconsole
96 # %matplotlib qtconsole
97 # %matplotlib qtconsole
98 # %matplotlib qtconsole
99 # %matplotlib qtconsole
100 # %matplotlib qtconsole
```

Categorical data

A 7.5 magnitude earthquake on September 28, 2018 near Palu, Indonesia caused a devastating tsunami. Let's take a look at some visualizations to understand what `magType` are used in Indonesia, the range of magnitudes there, and how many of the earthquakes are accompanied by a tsunami.

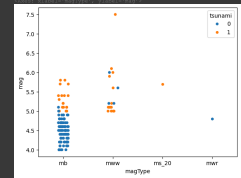
```
1 quakes.groupby('magType').size()
2 quakes.groupby('magType').size().reset_index()
3 quakes.groupby('magType').size().reset_index().sort_values('size', ascending=False)
4 quakes.groupby('magType').size().reset_index().sort_values('size', ascending=False).head(10)
```



stripplot()

The `stripplot()` function helps us visualize categorical data on one axis and numerical data on the other. We also now have the option of coloring our points using a column of our data (with the `hue` parameter). Using a strip plot, we can see points for each earthquake that was measured with a given `magType` and what its magnitude was. However, it isn't too easy to see density of the points due to overlap.

```
1 sns.stripplot(
2     data=quakes,
3     x='magType',
4     y='mag',
5     hue='tsunami',
6     jitter=True,
7     size=10,
8     style='white',
9     palette='magma',
10    )
```

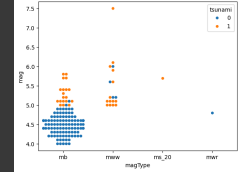


swarmplot()

The `swarm` plot helps address this issue by keeping the points from overlapping. Notice how many more points we can see for the blue section of the `mb` `magType`.

```
1 sns.swarmplot(
2     data=quakes,
3     x='magType',
4     y='mag',
5     hue='tsunami',
6     jitter=True,
7     size=10,
8     style='white',
9     palette='magma',
10    )
```

Notice that the `swarmplot()` function is similar to the `stripplot()` function, but it uses a different method to avoid overlap. The `swarmplot()` function is similar to the `stripplot()` function, but it uses a different method to avoid overlap. The `swarmplot()` function is similar to the `stripplot()` function, but it uses a different method to avoid overlap.

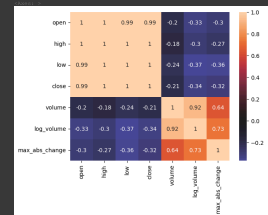


Correlations and Heatmaps

heatmap()

An easier way to create correlation matrix is to use `seaborn`:

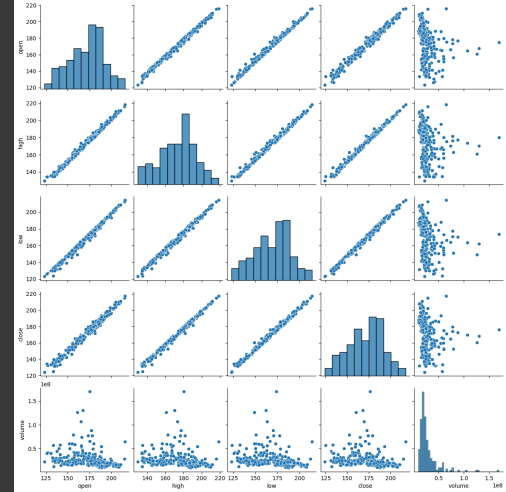
```
1 sns.heatmap(
2     data=quakes,
3     x='magType',
4     y='mag',
5     hue='tsunami',
6     jitter=True,
7     size=10,
8     style='white',
9     palette='magma',
10    )
```



pairplot()

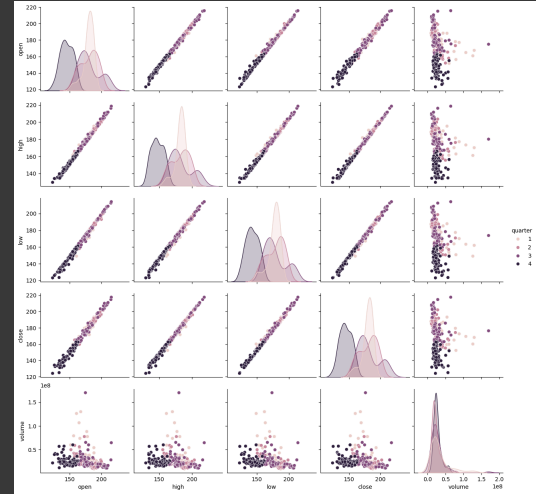
The pair plot is seaborn's answer to the scatter matrix we saw in the pandas subplotting notebook.

```
1 sns.pairplot(quakes)
```



Just as with pandas we can specify what to show along the diagonal; however, `seaborn` also allows us to color the data based on another column (or other data with the same shape).

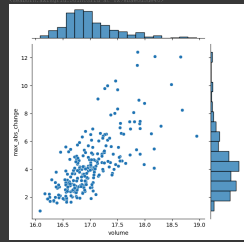
```
1 sns.pairplot(
2     data=quakes,
3     x='magType',
4     y='mag',
5     hue='tsunami',
6     jitter=True,
7     size=10,
8     style='white',
9     palette='magma',
10    )
```



jointplot()

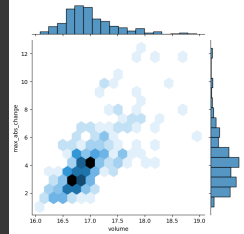
The joint plot allows us to visualize the relationship between two variables, like a scatter plot. However, we get the added benefit of being able to visualize their distributions at the same time (as a histogram or KDE). The default options give us a scatter plot in the center and histograms on the sides.

```
1 sns.jointplot(
2     data=quakes,
3     x='magType',
4     y='mag',
5     hue='tsunami',
6     jitter=True,
7     size=10,
8     style='white',
9     palette='magma',
10    )
```



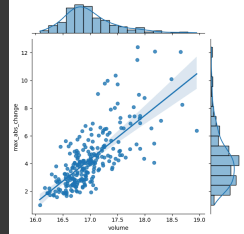
By changing the kind argument, we can change how the center of the plot is displayed. For example, we can pass kind='hex' for hexbins.

```
1 # Create a hexbin plot
2 fig, ax = plt.subplots()
3 ax.hist2d(volume, max_abs_change,
4           bins=20,
5           range=[(16, 19), (2, 12)],
6           cmap='magma')
7 ax.set_xlabel('volume')
8 ax.set_ylabel('max_abs_change')
9 ax.set_title('reg_abs_change_hexbin')
```



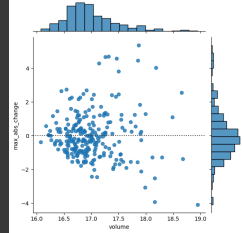
If we specify kind='reg', instead, we get a regression line in the center and KDEs on the sides.

```
1 # Create a regression plot
2 fig, ax = plt.subplots()
3 ax.hist2d(volume, max_abs_change,
4           bins=20,
5           range=[(16, 19), (2, 12)],
6           cmap='magma')
7 ax.set_xlabel('volume')
8 ax.set_ylabel('max_abs_change')
9 ax.set_title('reg_abs_change_reg')
```



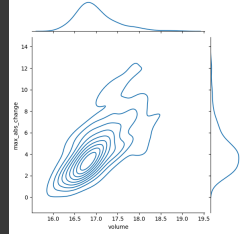
If we pass kind='resid', we get the residuals from the aforementioned regression.

```
1 # Create a residual plot
2 fig, ax = plt.subplots()
3 ax.hist2d(volume, max_abs_change,
4           bins=20,
5           range=[(16, 19), (2, 12)],
6           cmap='magma')
7 ax.set_xlabel('volume')
8 ax.set_ylabel('max_abs_change')
9 ax.set_title('reg_abs_change_resid')
```



Finally, if we pass kind='kde', we get a contour plot of the joint density estimate with KDEs along the sides.

```
1 # Create a KDE plot
2 fig, ax = plt.subplots()
3 ax.hist2d(volume, max_abs_change,
4           bins=20,
5           range=[(16, 19), (2, 12)],
6           cmap='magma')
7 ax.set_xlabel('volume')
8 ax.set_ylabel('max_abs_change')
9 ax.set_title('reg_abs_change_kde')
```



Regression plots

We are going to use `seaborn` to visualize a linear regression between the log of the volume traded in Facebook stock and the maximum absolute daily change (daily high stock price - daily low stock price). To do so, we first need to isolate this data.

```
1 # Import data
2 df = pd.read_csv('data/facebook_stock_data.csv')
3 df = df[df['volume'] > 10]
4 df = df[df['max_abs_change'] > 0]
```

Since we want to visualize each column as the regressor, we need to look at permutations of their order. Permutations and combinations (among other things) are made easy in Python with `itertools`, so let's import it:

```
1 import itertools
```

`itertools` gives us efficient iterators. Iterators are objects that we loop over, exhausting them. This is an iterator from `itertools`, notice how the second loop doesn't do anything:

```
1 iterator = itertools.repeat('it's an iterator', 5)
2 for i in iterator:
3     print(i)
4 # prints('it's an iterator')
5 for i in iterator:
6     print(i)
7 # prints('')
8 # prints('')
9 # prints('')
```

Iterables are objects that can be iterated over. When entering a loop, an iterator is made from the iterable to handle the iteration. Iterators are *iterables*, but not all iterables are iterators. A list is an *iterable*. If we turn that iterator into an *iterable* (a list in this case), the second loop runs:

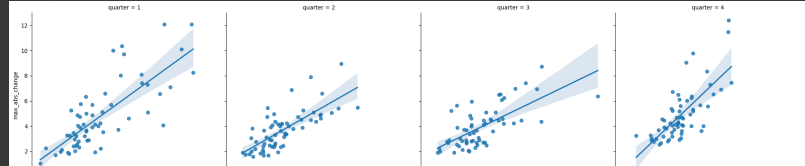
```
1 iterator = itertools.repeat('it's an iterator', 5)
2 for i in iterator:
3     print(i)
4 # prints('it's an iterator')
5 for i in iterator:
6     print(i)
7 # prints('it's an iterator')
8 # prints('it's an iterator')
9 # prints('it's an iterator')
```

The `reg_resid_plot()` function from the `reg_resid_plot.py` module in this folder uses `regplot()` and `residplot()` from `seaborn` along with `itertools` to plot the regression and residuals side-by-side.

```
1 # Import data
2 df = pd.read_csv('data/facebook_stock_data.csv')
3 df = df[df['volume'] > 10]
4 df = df[df['max_abs_change'] > 0]
5 # Create a regression plot
6 fig, ax = plt.subplots()
7 ax.hist2d(volume, max_abs_change,
8           bins=20,
9           range=[(16, 19), (2, 12)],
10          cmap='magma')
11 ax.set_xlabel('volume')
12 ax.set_ylabel('max_abs_change')
13 ax.set_title('reg_abs_change_reg')
```

We can use `lmplot()` to split our regression across subsets of our data. For example, we can perform a regression per quarter on the Facebook stock data.

```
1 # Import data
2 df = pd.read_csv('data/facebook_stock_data.csv')
3 df = df[df['volume'] > 10]
4 df = df[df['max_abs_change'] > 0]
5 # Create a regression plot
6 fig, ax = plt.subplots()
7 ax.hist2d(volume, max_abs_change,
8           bins=20,
9           range=[(16, 19), (2, 12)],
10          cmap='magma')
11 ax.set_xlabel('volume')
12 ax.set_ylabel('max_abs_change')
13 ax.set_title('reg_abs_change_reg')
```

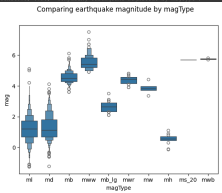


Distributions

`Seaborn` provides some new plot types for visualizing distributions in addition to its own versions of the plot types we discussed in chapter 5 (in this notebook).

The boxplot is a box plot that shows additional quantiles

```
1 ans.boxplot()
2
3 # magType = "mag", dataquest(magType, "mag")
4
5 plt.gca().set_title("Comparing earthquake magnitude by magType")
```

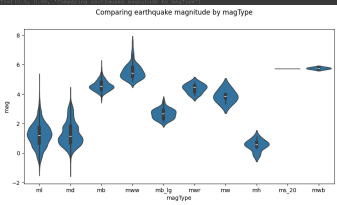


violinplot()

Box plots lose some information about the distribution, so we can use violin plots which combine box plots and K

```
1 fig, ax = plt.subplots(figsize=(10, 5))
2 ax.violinplot()
3
4 # magType = "mag", dataquest(magType, "mag")
5
6 # ax = ax, title = "title" # All violin have same width
7
8
9 plt.gca().set_title("Comparing earthquake magnitude by magType")
```

The "width" parameter has been removed and will be removed in 0.15.0. Pass "density_norm='width'" for the same effect.



Faceting

We can create subplots across subsets of our data by faceting. First, we create a FacetGrid specifying how to layout the plots (which categorical column goes along the rows and which one along the columns). Then, we call the map() method of the FacetGrid and pass in the plotting function we want to use (along with any additional arguments).

Let's make histograms showing the distribution of earthquake magnitude in California, Alaska, and Hawaii faceted by magType and parse_place:

```
1 g = sns.FacetGrid()
2 g.set()
3
4 # parse_place = "place", dataquest(place, "place")
5
6 # California, "Alaska", "Hawaii"
7
8 # 10
9
10 # parse_place = "place", dataquest(place, "place")
11
12 # 10
13
14 # 10
15
16 # 10
17
18 # 10
19
20 # 10
21
22 # 10
23
24 # 10
25
26 # 10
27
28 # 10
29
30 # 10
31
32 # 10
33
34 # 10
35
36 # 10
37
38 # 10
39
40 # 10
41
42 # 10
43
44 # 10
45
46 # 10
47
48 # 10
49
50 # 10
51
52 # 10
53
54 # 10
55
56 # 10
57
58 # 10
59
60 # 10
61
62 # 10
63
64 # 10
65
66 # 10
67
68 # 10
69
70 # 10
71
72 # 10
73
74 # 10
75
76 # 10
77
78 # 10
79
80 # 10
81
82 # 10
83
84 # 10
85
86 # 10
87
88 # 10
89
90 # 10
91
92 # 10
93
94 # 10
95
96 # 10
97
98 # 10
99
100 # 10
```

