

9.4 Introduction to Seaborn

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

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import pandas as pd
fb = pd.read_csv(
    'fb_stock_prices_2018.csv', index_col='date', parse_dates=True
)
quakes = pd.read_csv('earthquakes.csv')
```


Categorical data

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

```
quakes.assign(
    time=lambda x: pd.to_datetime(x.time, unit='ms')
).set_index('time').loc['2018-09-28'].query(
    "parsed_place == 'Indonesia' and tsunami == 1 and mag == 7.5"
)
```



	mag	magType	place	tsunami	parsed_place
time					
2018-09-28 10:02:43.480	7.5	mww	78km N of Palu, Indonesia	1	Indonesia

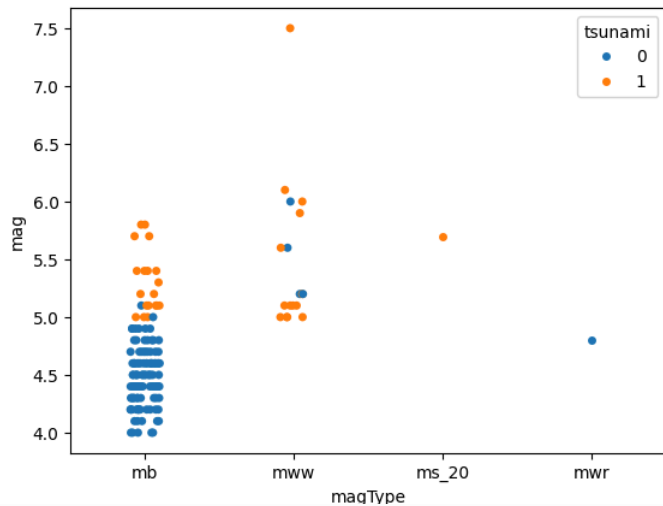


stripplot()

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 was; however, it isn't too easy to see density of the points due to overlap

```
sns.stripplot(
    x='magType',
    y='mag',
    hue='tsunami',
    data=quakes.query('parsed_place == "Indonesia"')
)
```

```
>>> <Axes: xlabel='magType', ylabel='mag'>
```



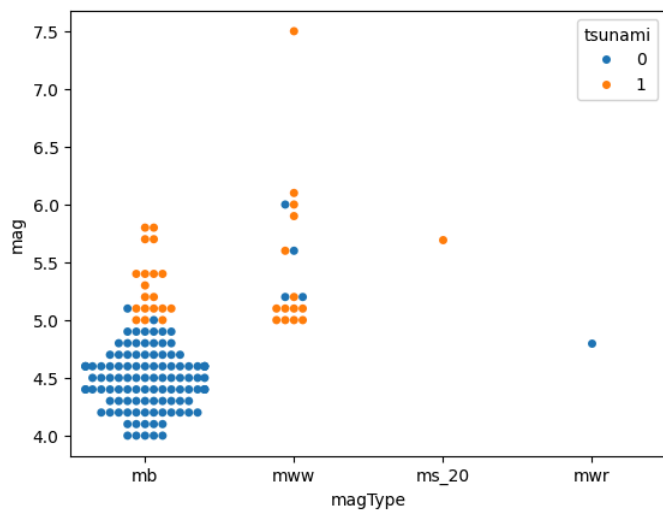
swarmplot

The bee 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 :

```
sns.swarmplot(
    x='magType',
    y='mag',
    hue='tsunami',
    data=quakes.query('parsed_place == "Indonesia"')
)
```

```
>>> <Axes: xlabel='magType', ylabel='mag'>
```

/usr/local/lib/python3.11/dist-packages/seaborn/categorical.py:3399: UserWarning: 10.2% of the points cannot be placed; you may want to decrease the size of the markers or the number of points.
warnings.warn(msg, UserWarning)



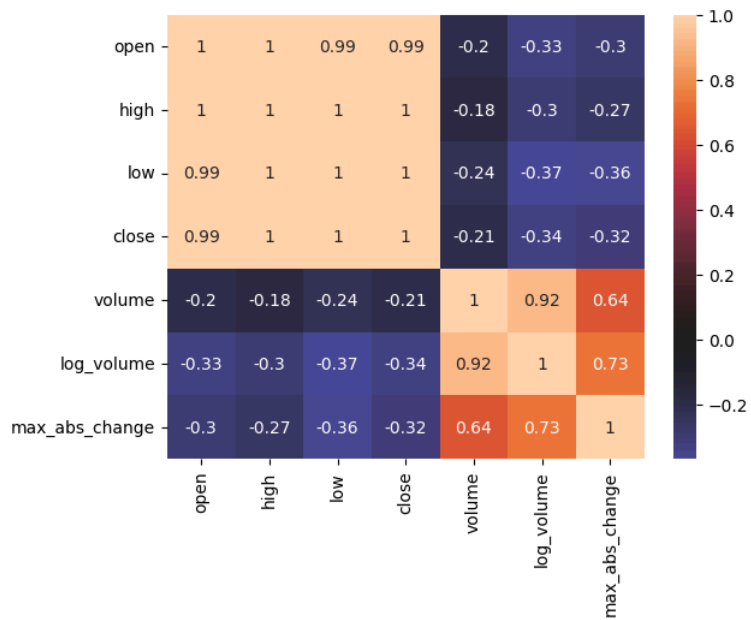
Correlations and Heatmaps

heatmap()

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

```
sns.heatmap(
    fb.sort_index().assign(
        log_volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low
    ).corr(),
    annot=True, center=0
)
```

<Axes: >

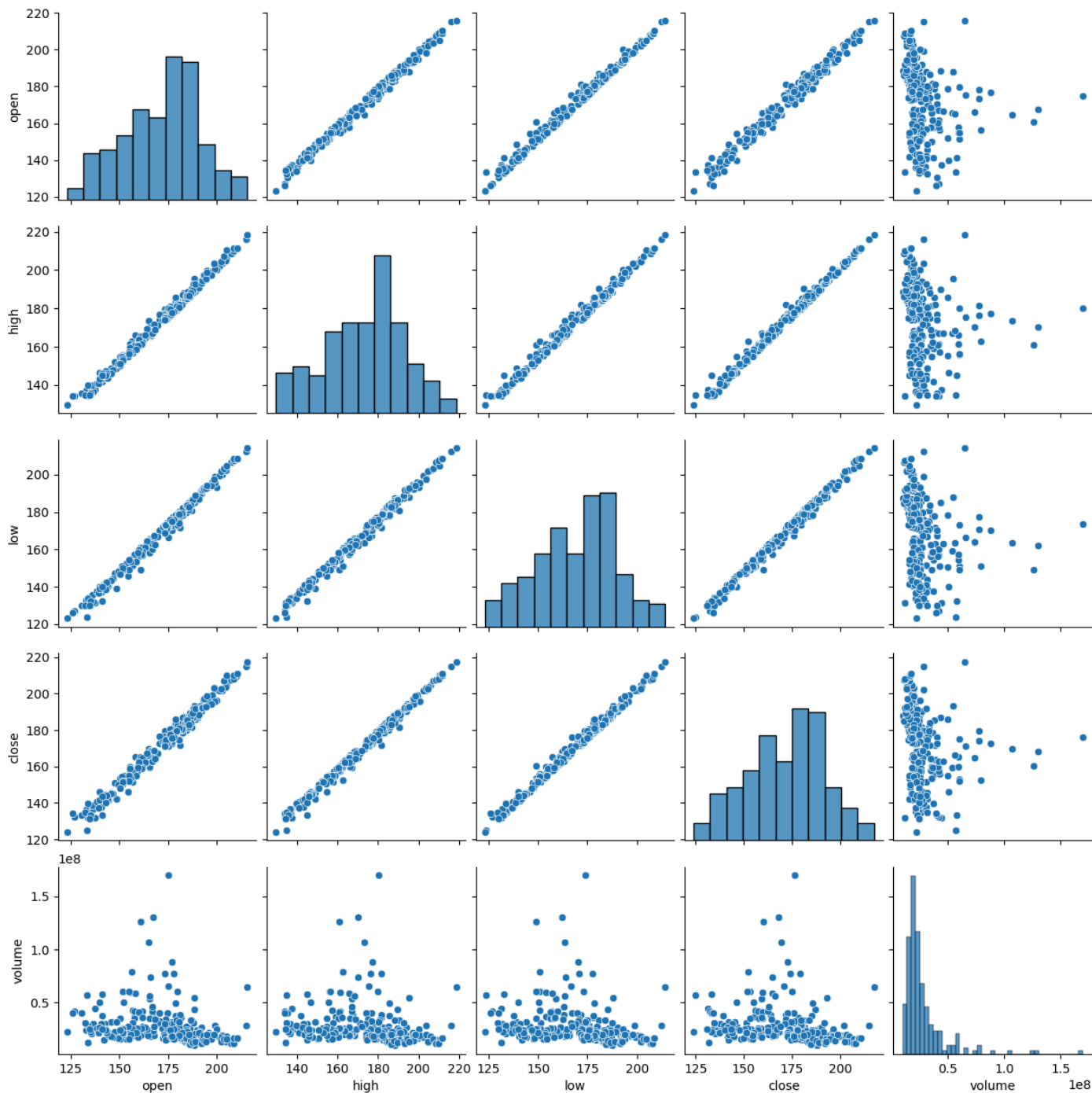


pairplot()

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

sns.pairplot(fb)

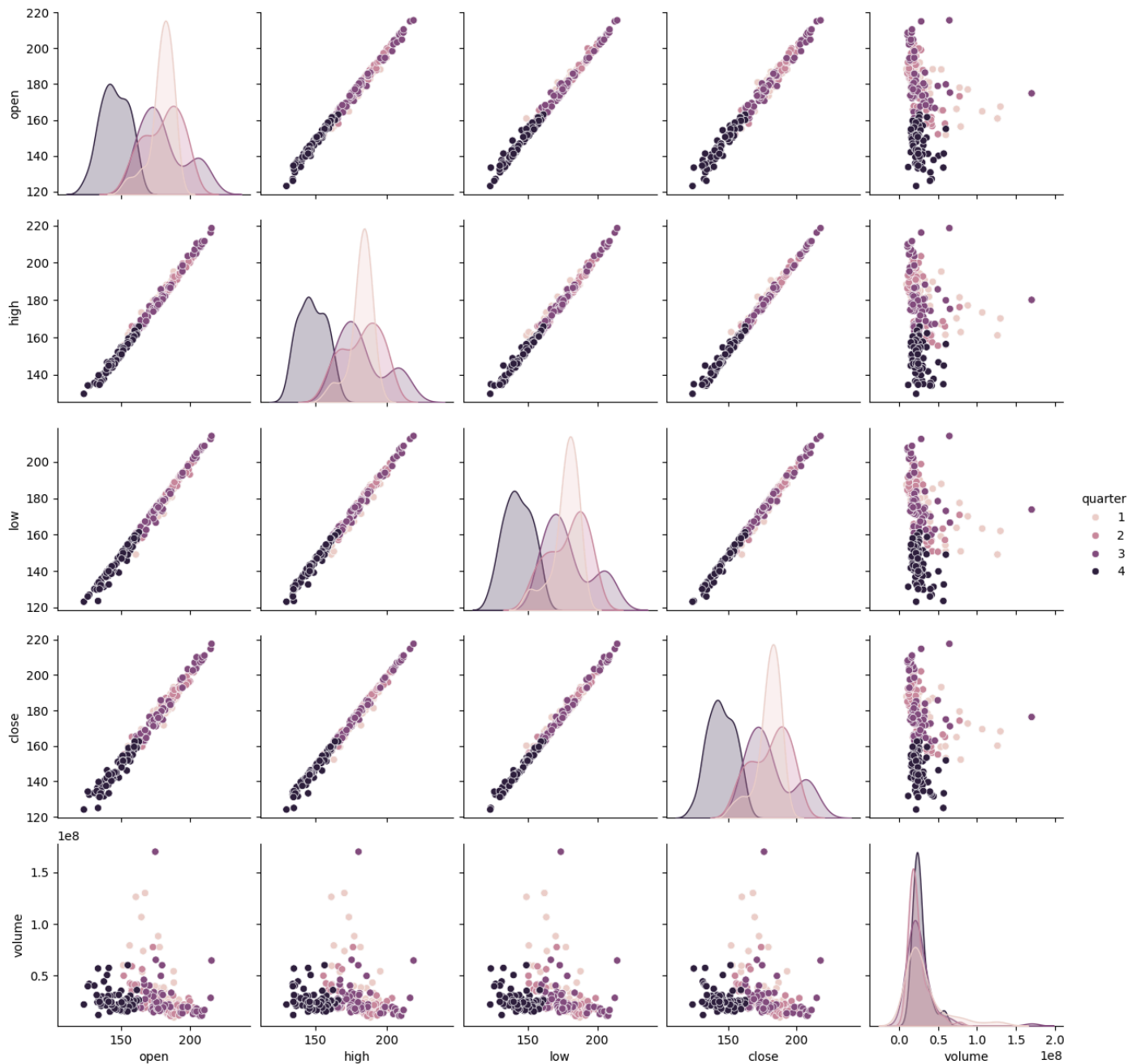
`<seaborn.axisgrid.PairGrid at 0x7d20f3b7bed0>`



Just as with same shape): 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

```
sns.pairplot(
    fb.assign(quarter=lambda x: x.index.quarter),
    diag_kind='kde',
    hue='quarter'
)
```

```
<seaborn.axisgrid.PairGrid at 0x7d20f3f941d0>
```

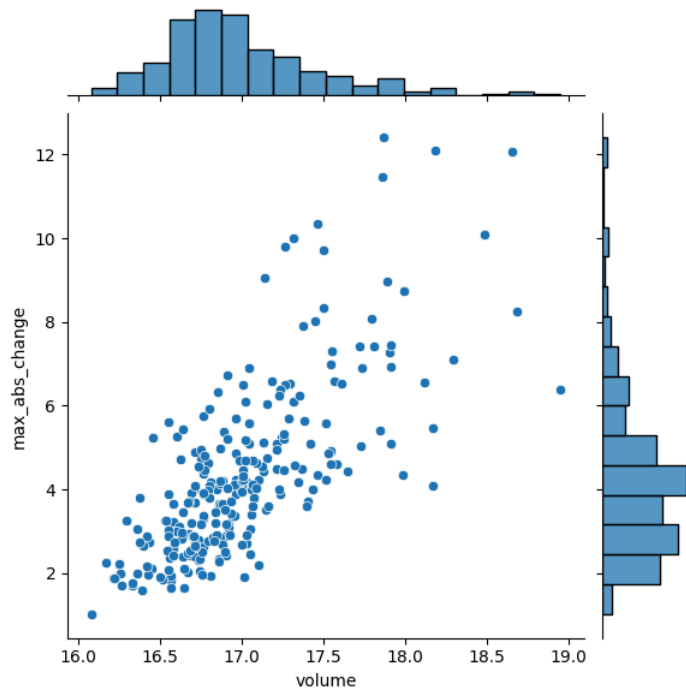


```
jointplot()
```

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:

```
sns.jointplot(
    x='volume',
    y='max_abs_change',
    data=fb.assign(
        volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low
    )
)
```

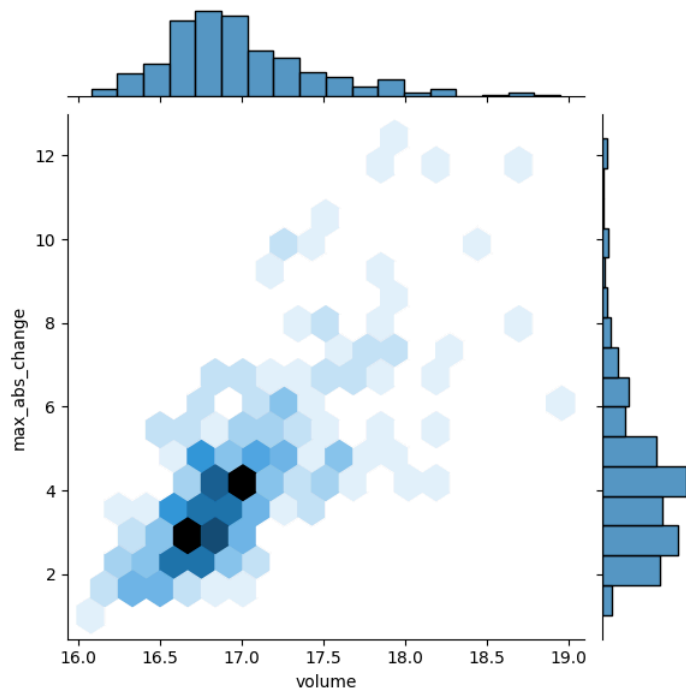
```
>>> sns.jointplot(x='volume', y='max_abs_change', kind='scatter', data=fb, size=10)
```



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:

```
sns.jointplot(
    x='volume',
    y='max_abs_change',
    kind='hex',
    data=fb.assign(
        volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low
    )
)
```

```
>>> sns.jointplot(x='volume', y='max_abs_change', kind='hex', data=fb, size=10)
```



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

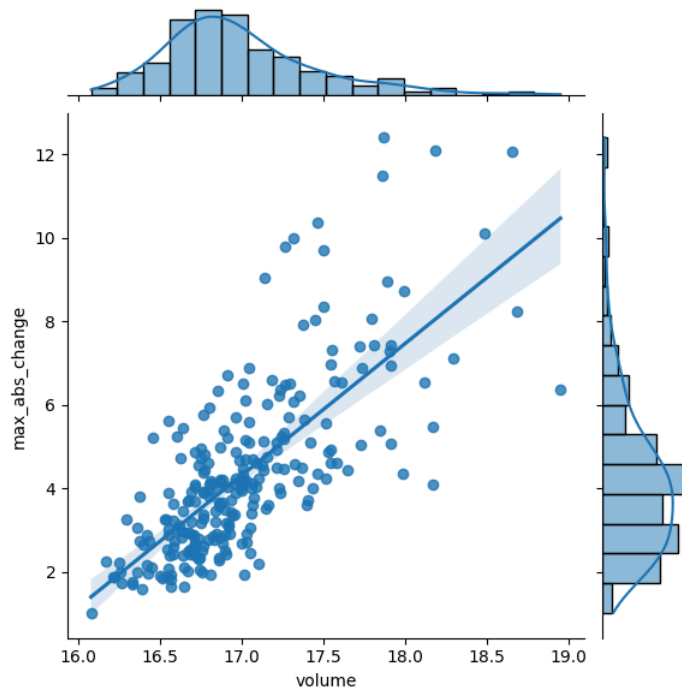
```
sns.jointplot(
    x='volume',
    y='max_abs_change',
    kind='reg',
    data=fb,
    size=10
)
```

```

y='max_abs_change',
kind='reg',
data=fb.assign(
    volume=np.log(fb.volume),
    max_abs_change=fb.high - fb.low
)
)

```

 <seaborn.axisgrid.JointGrid at 0x7d20f2bca810>



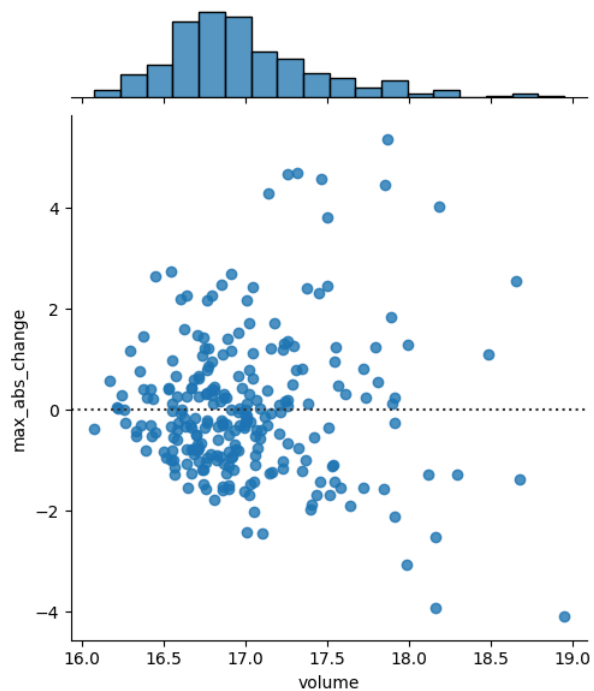
If we pass kind='resid', we get the residuals from the aforementioned regression

```

sns.jointplot(
    x='volume',
    y='max_abs_change',
    kind='resid',
    data=fb.assign(
        volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low
    )
)

```

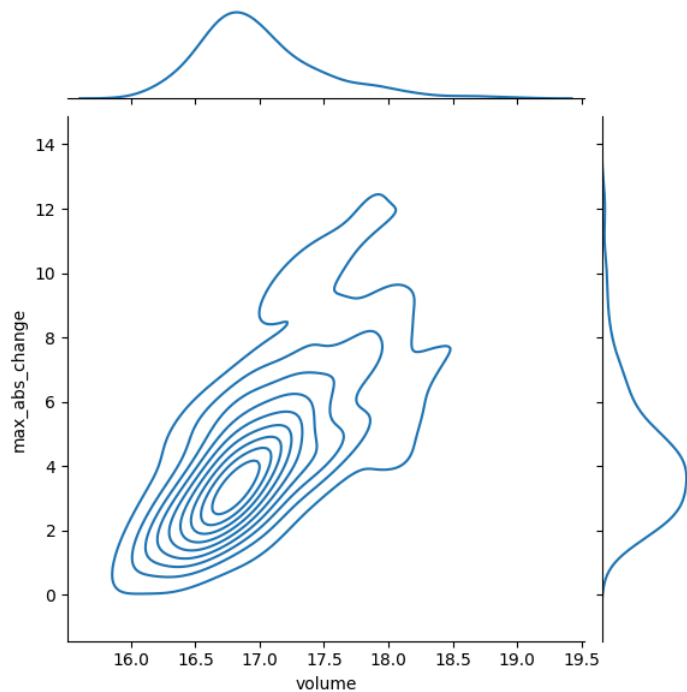
```
<seaborn.axisgrid.JointGrid at 0x7d20f2b26290>
```



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

```
sns.jointplot(
    x='volume',
    y='max_abs_change',
    kind='kde',
    data=fb.assign(
        volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low
    )
)
```

```
<seaborn.axisgrid.JointGrid at 0x7d20f2915b90>
```



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


```
fb_reg_data = fb.assign(
    volume=np.log(fb.volume),
    max_abs_change=fb.high - fb.low
).iloc[:, -2:]
```

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

```
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:

```
iterator = itertools.repeat("I'm an iterator", 1)
for i in iterator:
    print(f'-->{i}')
print('This printed once because the iterator has been exhausted')
for i in iterator:
    print(f'-->{i}')
```

```
-->I'm an iterator
This printed once because the iterator has been exhausted
```

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:

```
iterable = list(itertools.repeat("I'm an iterable", 1))
for i in iterable:
    print(f'-->{i}')
print('This prints again because it's an iterable:')
for i in iterable:
    print(f'-->{i}')
```

```
-->I'm an iterable
This prints again because it's an iterable:
-->I'm an iterable
```

The `reg_resid_plots()` 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

```
from reg_resid_plot import reg_resid_plots
reg_resid_plots(fb_reg_data)
```

```
-----
ModuleNotFoundError                                Traceback (most recent call last)
<ipython-input-155-ae2d095ec697> in <cell line: 0>()
----> 1 from reg_resid_plot import reg_resid_plots
      2 reg_resid_plots(fb_reg_data)
```

```
ModuleNotFoundError: No module named 'reg_resid_plot'
```

```
-----
NOTE: If your import is failing due to a missing package, you can
manually install dependencies using either !pip or !apt.
```

```
To view examples of installing some common dependencies, click the
"Open Examples" button below.
```

OPEN EXAMPLES

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

```
sns.lmplot(
    x='volume',
    y='max_abs_change',
    data=fb.assign(
        volume=np.log(fb.volume),
        max_abs_change=fb.high - fb.low,
        quarter=lambda x: x.index.quarter
    ),
    col='quarter'
)
```

```

AttributeError                                Traceback (most recent call last)
<ipython-input-157-3f2bf7175838> in <cell line: 0>()
      2 x='volume',
      3 y='max_abs_change',
----> 4 data=fb.assign(
      5     volume=np.log(fb.volume),
      6     max_abs_change=fb.high - fb.low,

2 frames
<ipython-input-157-3f2bf7175838> in <lambda>(x)
      5     volume=np.log(fb.volume),
      6     max_abs_change=fb.high - fb.low,
----> 7     quarter=lambda x: x.index.quarter
      8 ),
      9     col='quarter'

AttributeError: 'RangeIndex' object has no attribute 'quarter'

```

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)

boxenplot()

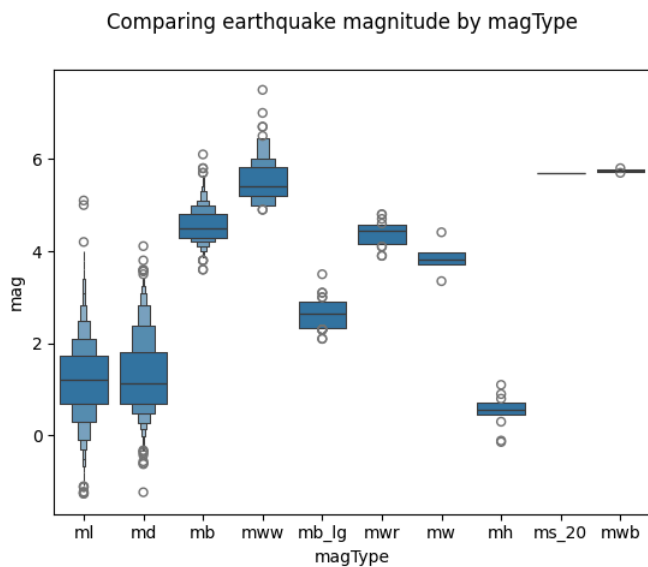
The boxenplot is a box plot that shows additional quantiles

```

sns.boxenplot(
    x='magType', y='mag', data=quakes[['magType', 'mag']]
)
plt.suptitle('Comparing earthquake magnitude by magType')

```

Text(0.5, 0.98, '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 KDEs

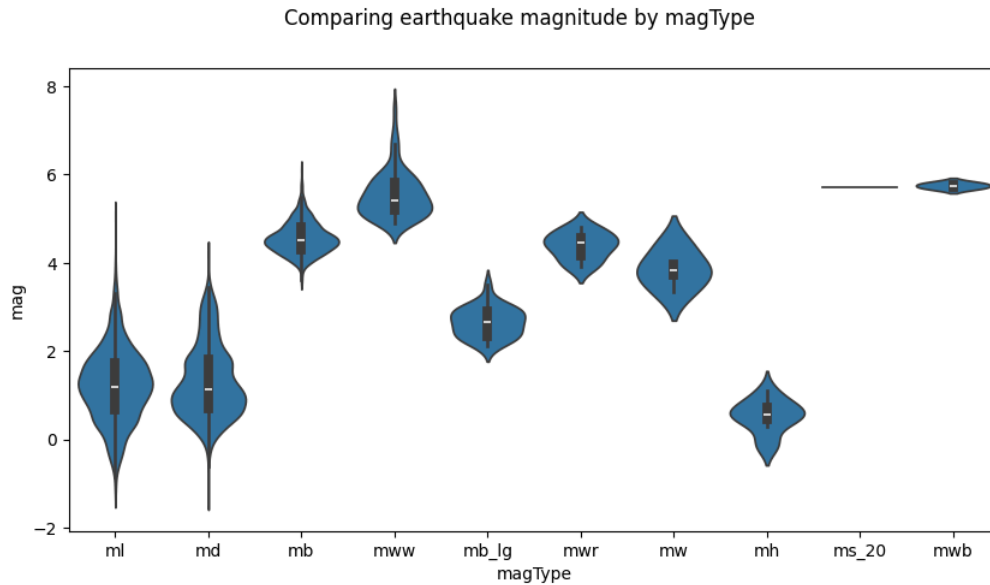
```

fig, axes = plt.subplots(figsize=(10, 5))
sns.violinplot(
    x='magType', y='mag', data=quakes[['magType', 'mag']],
    ax=axes, scale='width' # all violins have same width
)
plt.suptitle('Comparing earthquake magnitude by magType')

```

<ipython-input-159-774c9aedf3fb>:2: FutureWarning:

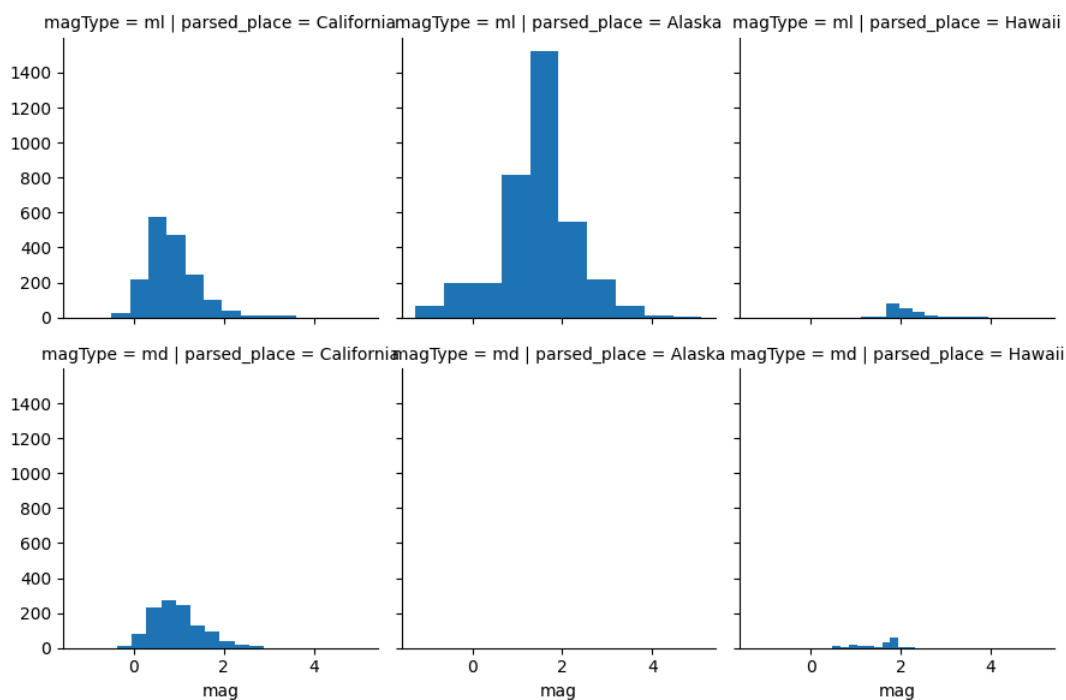
The `scale` parameter has been renamed and will be removed in v0.15.0. Pass `density_norm='width'` for the same effect.
sns.violinplot(
Text(0.5, 0.98, 'Comparing earthquake magnitude by magType')



Faceting

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

```
g = sns.FacetGrid(  
    quakes[  
        (quakes.parsed_place.isin([  
            'California', 'Alaska', 'Hawaii'  
        ]))\n        & (quakes.magType.isin(['ml', 'md']))  
    ],  
    row='magType',  
    col='parsed_place'  
)  
g = g.map(plt.hist, 'mag')
```



✓ 9.5 Formatting Plots

About the Data

In this notebook, we will be working with Facebook's stock price throughout 2018 (obtained using the `stock_analysis`

✓ Setup

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns

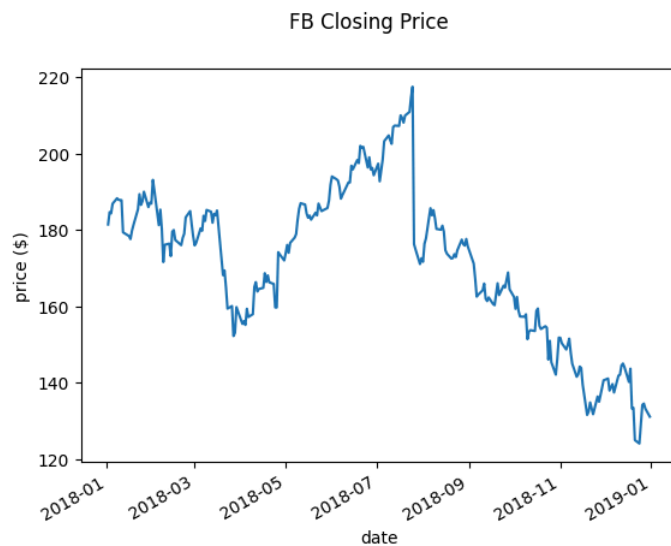
fb = pd.read_csv(
    'fb_stock_prices_2018.csv', index_col='date', parse_dates=True
)
```

✓ Titles and Axis Labels

`plt.suptitle()` adds a title to plots and subplots `plt.title()` adds a title to a single plot. Note if you use subplots, it will only put the title on the last subplot, so you will need to use `plt.xlabel()` labels the x-axis `plt.ylabel()` labels the y-axis

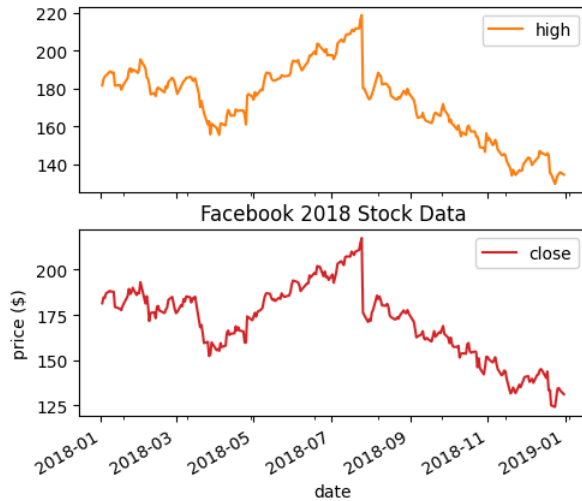
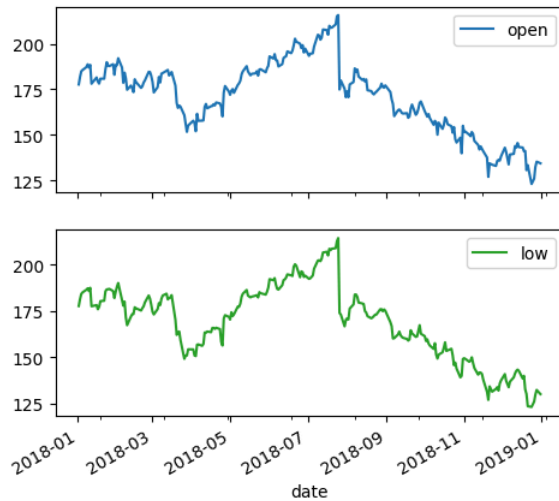
```
fb.close.plot()
plt.suptitle('FB Closing Price')
plt.xlabel('date')
plt.ylabel('price ($)')
```

```
→ Text(0, 0.5, 'price ($)')
```

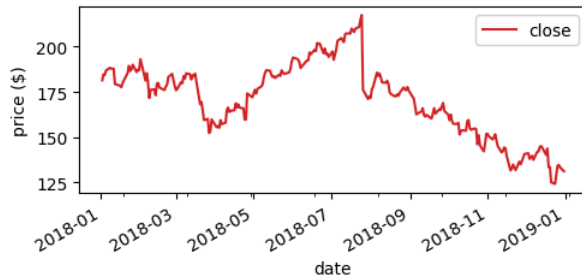
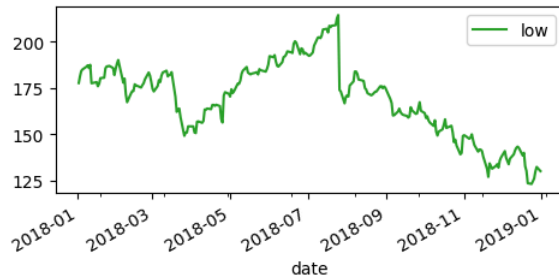


```
# plt.title
fb.iloc[:, :4].plot(subplots=True, layout=(2, 2), figsize=(12, 5))
plt.title('Facebook 2018 Stock Data')
plt.xlabel('date')
plt.ylabel('price ($)')
```

```
Text(0, 0.5, 'price ($)')
```



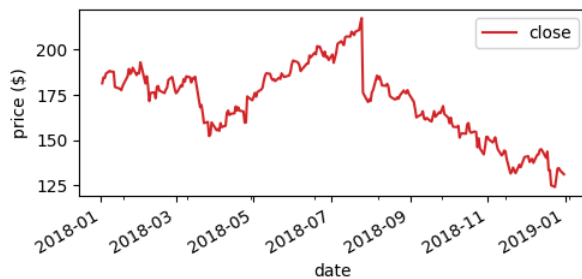
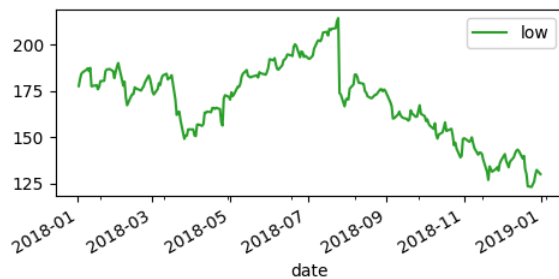
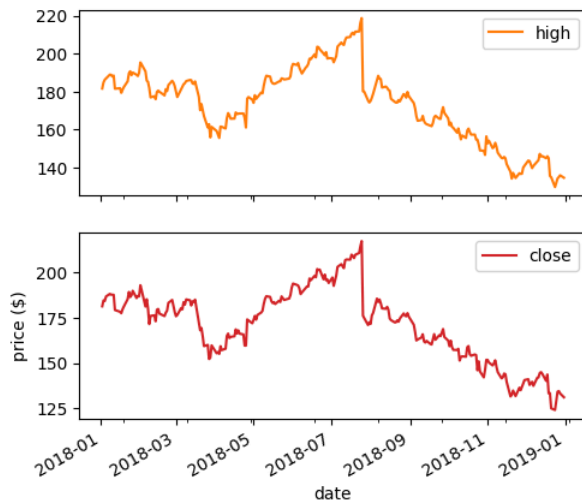
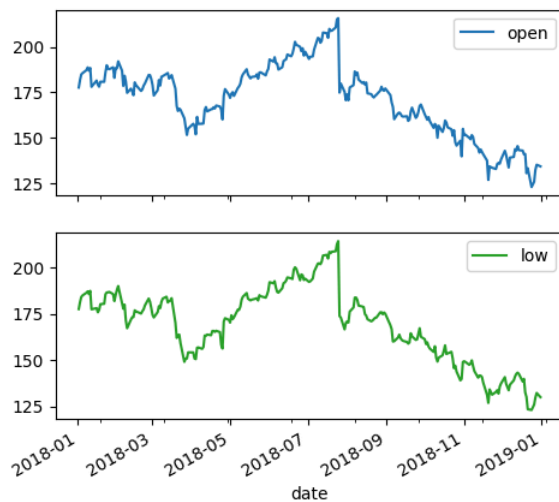
Facebook 2018 Stock Data



```
# plt.suptitle
fb.iloc[:,4].plot(subplots=True, layout=(2, 2), figsize=(12, 5))
plt.suptitle('Facebook 2018 Stock Data')
plt.xlabel('date')
plt.ylabel('price ($)')
```

```
Text(0, 0.5, 'price ($)')
```

Facebook 2018 Stock Data

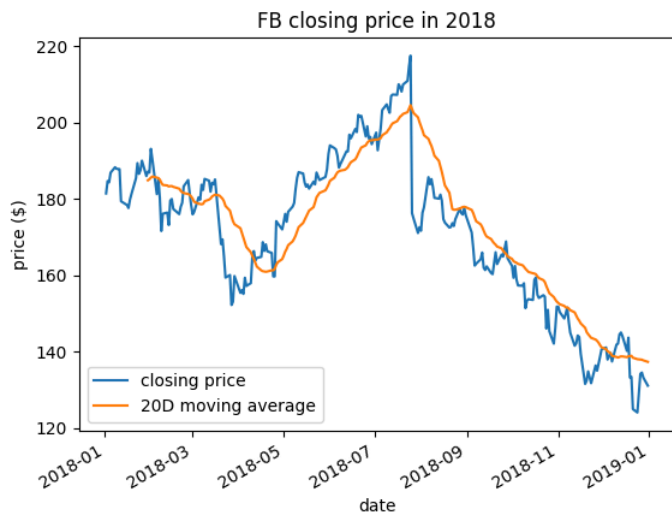


Legends

`plt.legend` adds a legend to the plot. We can specify where to place it with the `loc` parameter:

```
fb.assign(
    ma=lambda x: x.close.rolling(20).mean()
).plot(
    y=['close', 'ma'],
    title='FB closing price in 2018',
    label=['closing price', '20D moving average']
)
plt.legend(loc='lower left')
plt.ylabel('price ($)')
```

```
Text(0, 0.5, 'price ($)')
```

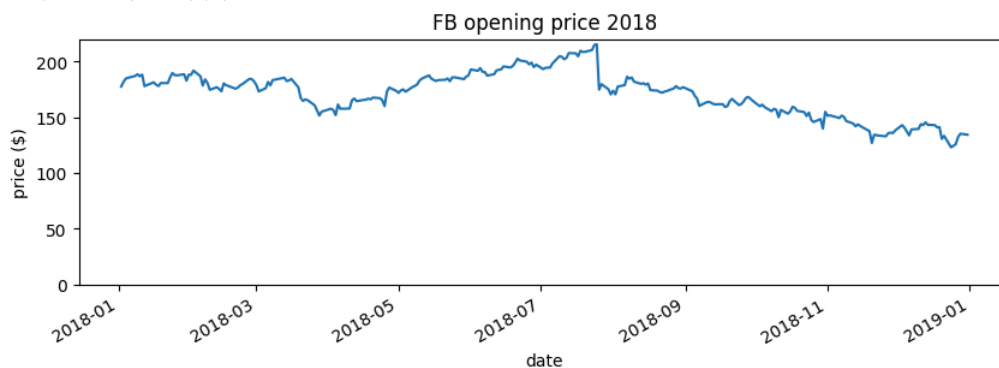


Formatting Axes

Specifying axis limits `plt.xlim()` and `plt.ylim()` can be used to specify the minimum and maximum values for the axis. Passing `None` will have matplotlib determine the limit

```
fb.open.plot(figsize=(10, 3), title='FB opening price 2018')
plt.ylim(0, None)
plt.ylabel('price ($)')
```

```
Text(0, 0.5, 'price ($)')
```



Formatting the axis Ticks

We can use `plt.xticks()` and `plt.yticks()` to provide tick labels and specify, which ticks to show. Here, we show every other month:

```
import calendar
fb.open.plot(figsize=(10, 3), rot=0, title='FB opening price 2018')
locs, labels = plt.xticks()
plt.xticks(locs + 15, calendar.month_name[1::2])
plt.ylabel('price ($)')
```

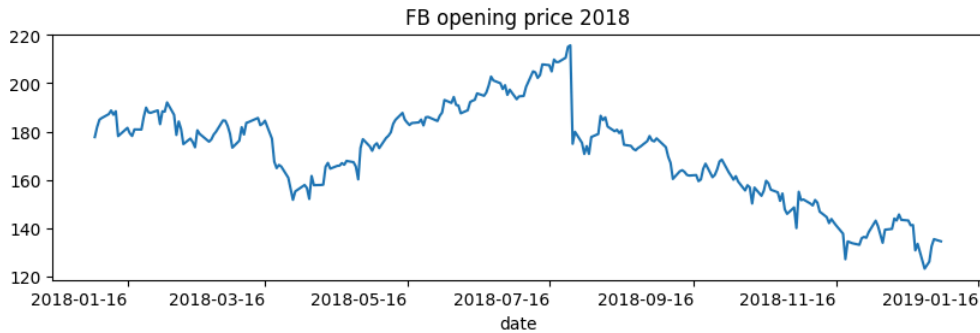
```

ValueError                                Traceback (most recent call last)
<ipython-input-166-49f9a03c7ca6> in <cell line: 0>()
      2 fb.open.plot(figsize=(10, 3), rot=0, title='FB opening price 2018')
      3 locs, labels = plt.xticks()
----> 4 plt.xticks(locs + 15, calendar.month_name[1::2])
      5 plt.ylabel('price ($)')

2 frames
/usr/local/lib/python3.11/dist-packages/matplotlib/axis.py in set_ticklabels(self, labels, minor, fontdict, **kwargs)
    2115         # remove all tick labels, so only error for > 0 labels
    2116         if len(locator.locs) != len(labels) and len(labels) != 0:
-> 2117             raise ValueError(
    2118                 "The number of FixedLocator locations"
    2119                 f" ({len(locator.locs)}), usually from a call to"

ValueError: The number of FixedLocator locations (7), usually from a call to set_ticks, does not match the number of labels (6).

```



```

import matplotlib.ticker as ticker

ax = fb.close.plot(
    figsize=(10, 4),
    title='Facebook Closing Price as Percentage of Highest Price in Time Range'
)
ax.yaxis.set_major_formatter(
    ticker.PercentFormatter(xmax=fb.high.max())
)
ax.set_yticks([
    fb.high.max()*pct for pct in np.linspace(0.6, 1, num=5)
]) # show round percentages only (60%, 80%, etc.)
ax.set_ylabel(f'percent of highest price (${fb.high.max()})')

Text(0, 0.5, 'percent of highest price ($218.62)')

```



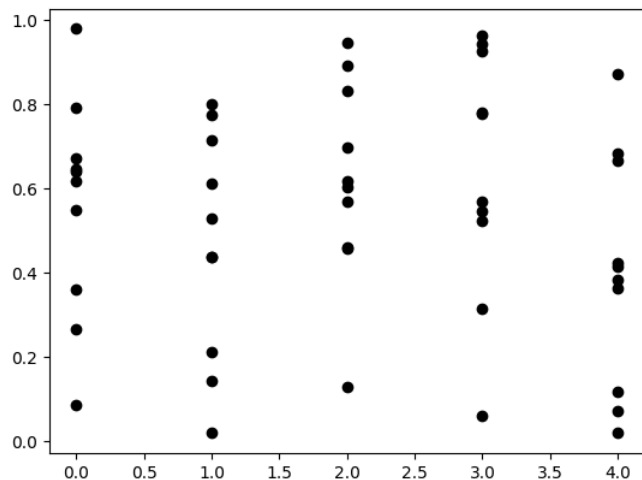
MultipleLocator

```

fig, ax = plt.subplots(1, 1)
np.random.seed(0)
ax.plot(np.tile(np.arange(0, 5), 10), np.random.rand(50), 'ko')

```

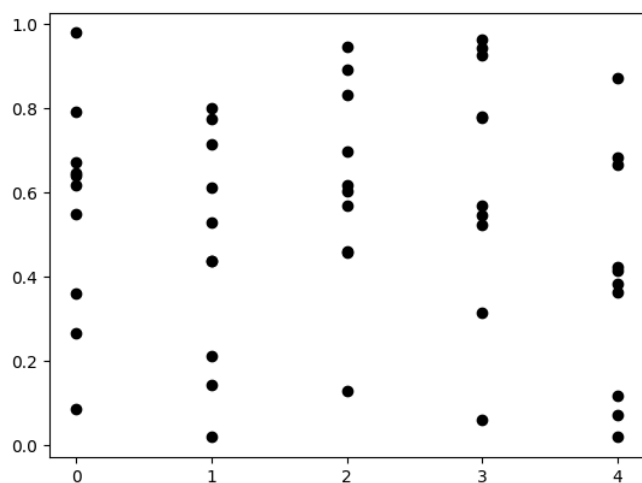
→ [<matplotlib.lines.Line2D at 0x7d20f1e50750>]



If we don't want to show decimal values on the x-axis, we can use the parameter. To get integer values, we use base=1 :

```
fig, ax = plt.subplots(1, 1)
np.random.seed(0)
ax.plot(np.tile(np.arange(0, 5), 10), np.random.rand(50), 'ko')
ax.get_xaxis().set_major_locator(
    ticker.MultipleLocator(base=1)
)
```

→



✓ 9.6 pandas.plotting subpackage

✓ Setup

```
%matplotlib inline
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd

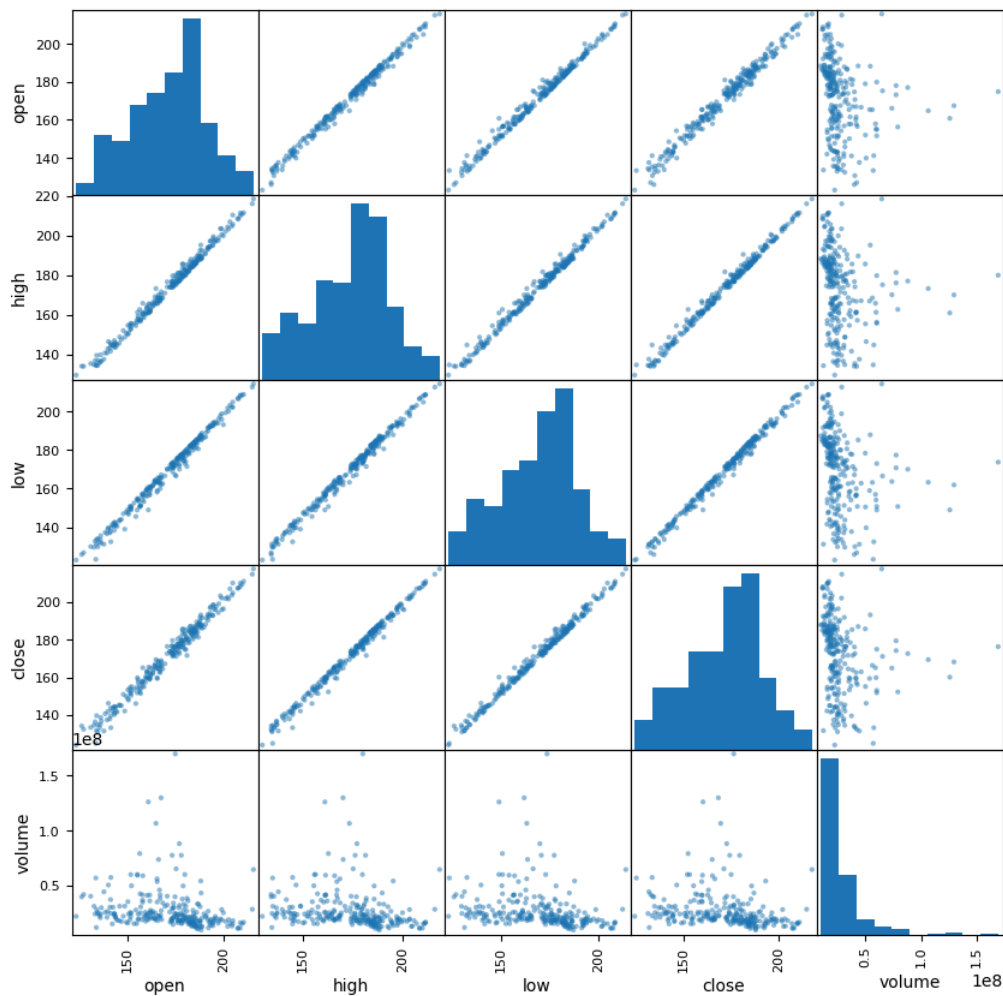
fb = pd.read_csv(
    'fb_stock_prices_2018.csv', index_col='date', parse_dates=True)
```

✓ Scatter Matri

```
from pandas.plotting import scatter_matrix
scatter_matrix(fb, figsize=(10, 10))
```



```
array([[<Axes: xlabel='open', ylabel='open'>,
        <Axes: xlabel='high', ylabel='open'>,
        <Axes: xlabel='low', ylabel='open'>,
        <Axes: xlabel='close', ylabel='open'>,
        <Axes: xlabel='volume', ylabel='open'>],
       [<Axes: xlabel='open', ylabel='high'>,
        <Axes: xlabel='high', ylabel='high'>,
        <Axes: xlabel='low', ylabel='high'>,
        <Axes: xlabel='close', ylabel='high'>,
        <Axes: xlabel='volume', ylabel='high'>],
       [<Axes: xlabel='open', ylabel='low'>,
        <Axes: xlabel='high', ylabel='low'>,
        <Axes: xlabel='low', ylabel='low'>,
        <Axes: xlabel='close', ylabel='low'>,
        <Axes: xlabel='volume', ylabel='low'>],
       [<Axes: xlabel='open', ylabel='close'>,
        <Axes: xlabel='high', ylabel='close'>,
        <Axes: xlabel='low', ylabel='close'>,
        <Axes: xlabel='close', ylabel='close'>,
        <Axes: xlabel='volume', ylabel='close'>],
       [<Axes: xlabel='open', ylabel='volume'>,
        <Axes: xlabel='high', ylabel='volume'>,
        <Axes: xlabel='low', ylabel='volume'>,
        <Axes: xlabel='close', ylabel='volume'>,
        <Axes: xlabel='volume', ylabel='volume'>]], dtype=object)
```




- ▼ Lag plot

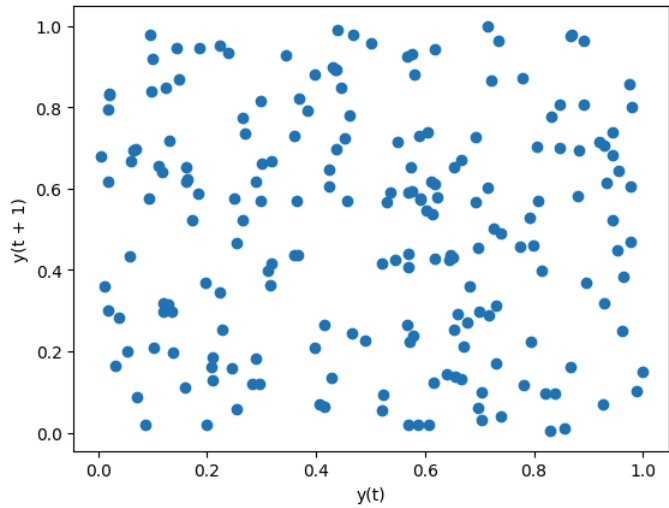
Lag plots let us see how the variable correlations with past observations of itself. Random data has no pattern:

```
from pandas.plotting import lag_plot
np.random.seed(0) # make this repeatable
lag_plot(pd.Series(np.random.random(size=200)))
```


```
np.random.seed(0) # make this repeatable
```

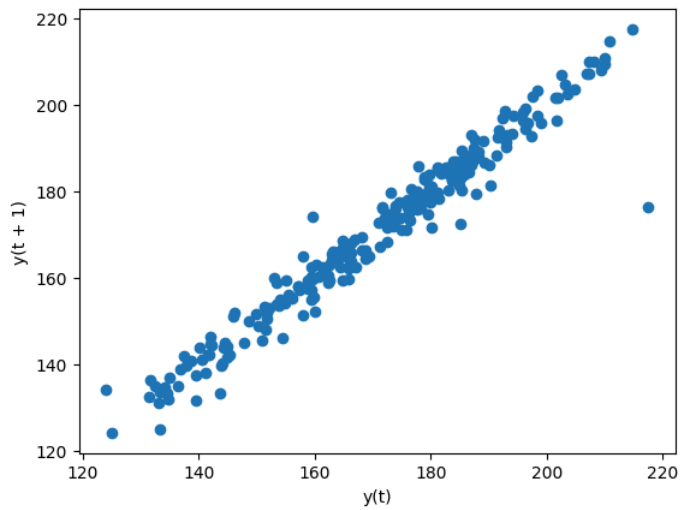
```
lag_plot(pd.Series(np.random.random(size=200)))
```

 <Axes: xlabel='y(t)', ylabel='y(t + 1)'\>




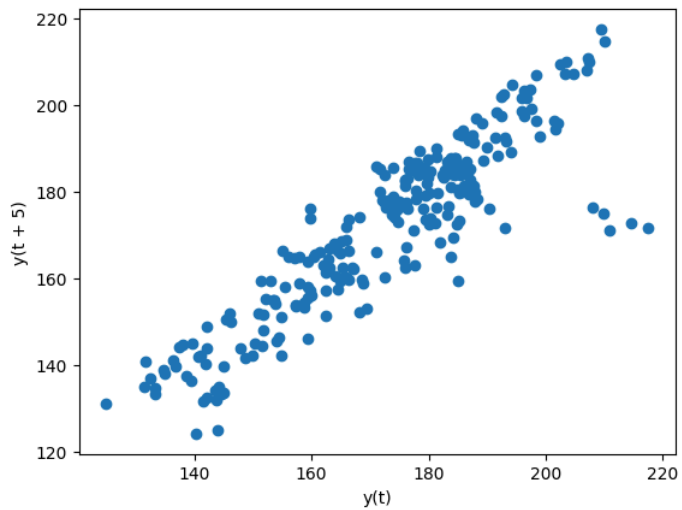
```
lag_plot(fb.close)
```

 <Axes: xlabel='y(t)', ylabel='y(t + 1)'\>



The default lag is 1, but we can alter this with the 'lag' parameter. Let's look at a 5 day lag (a week of trading activity) parameter. Let's look at a 5
`lag_plot(fb.close, lag=5)`

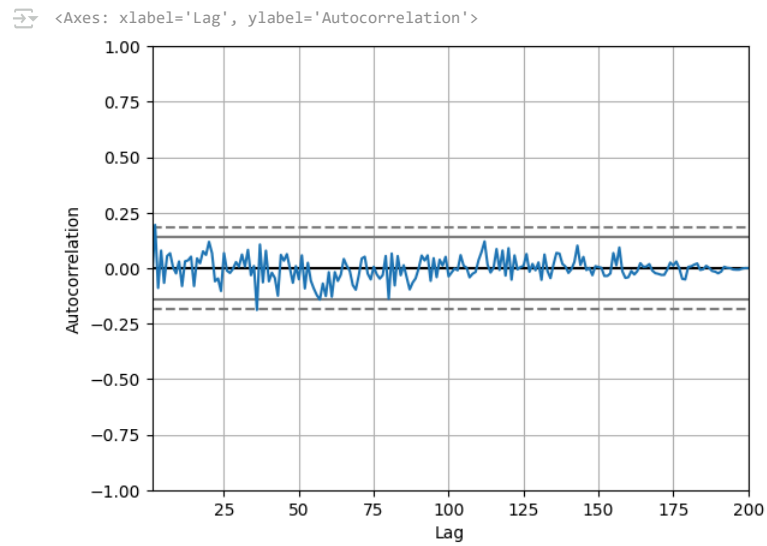
 <Axes: xlabel='y(t)', ylabel='y(t + 5)'\>



∨ Autocorrelation plots

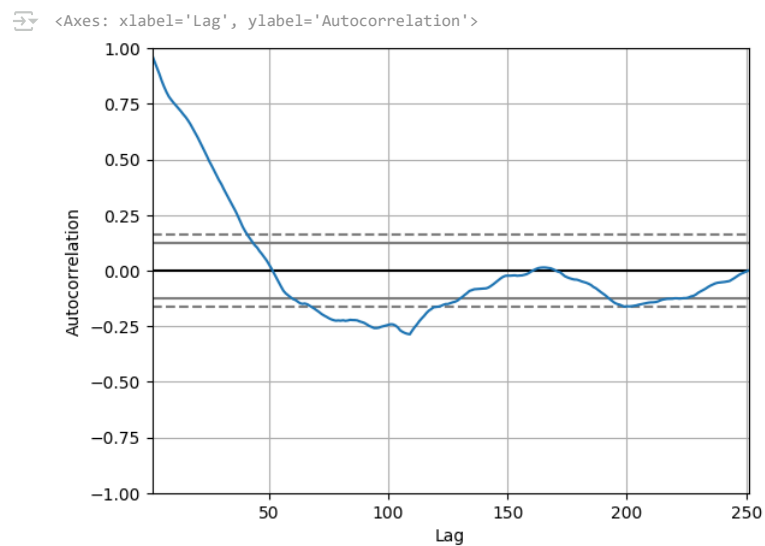
We can use the autocorrelation plot to see if this relationship may be meaningful or just noise. Random data will not have any significant autocorrelation (it stays within the bounds below):

```
from pandas.plotting import autocorrelation_plot
np.random.seed(0) # make this repeatable
autocorrelation_plot(pd.Series(np.random.random(size=200)))
```



Stock data, on the other hand, does have significant autocorrelation

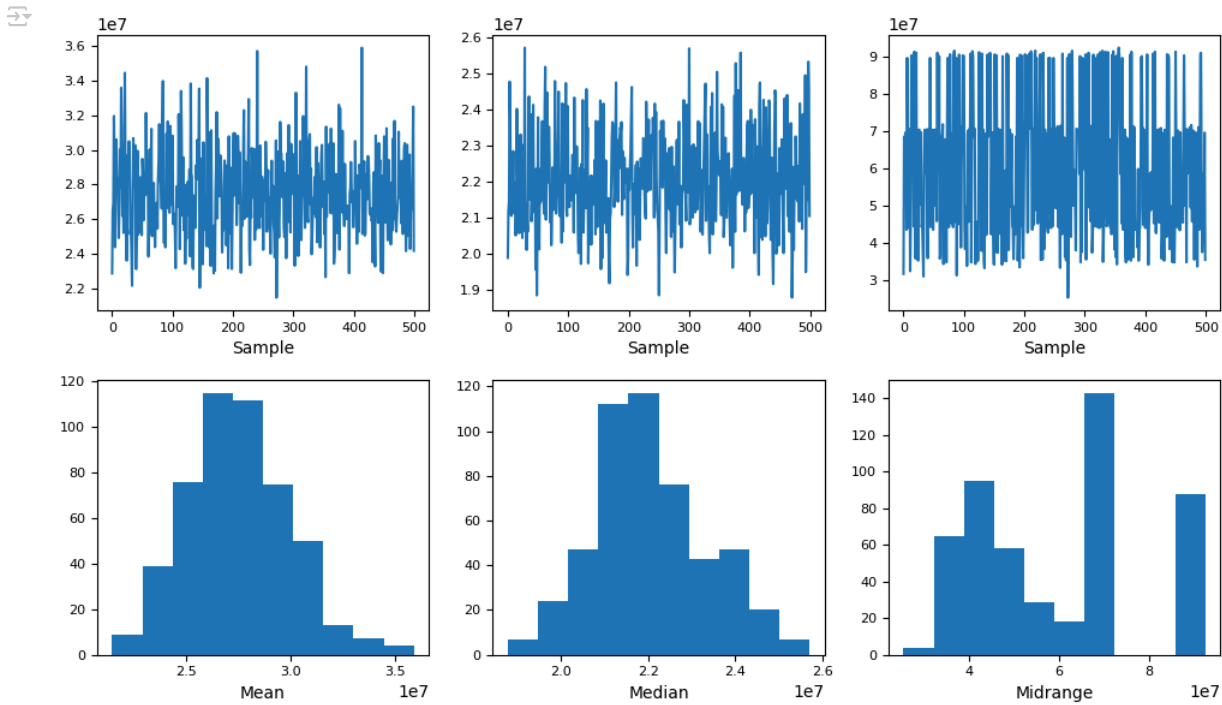
```
autocorrelation_plot(fb.close)
```



✓ Bootstrap plot

This plot helps us understand the uncertainty in our summary statistics:

```
from pandas.plotting import bootstrap_plot
fig = bootstrap_plot(fb.volume, fig=plt.figure(figsize=(10, 6)))
```



Supplementary Activity:

Using seaborn, create a heatmap to visualize the correlation coefficients between earthquake magnitude and whether there was a tsunami with the magType of mb.

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

quake = pd.read_csv('earthquakes.csv')
```

- 1.) Using seaborn, create a heatmap to visualize the correlation coefficients between earthquake magnitude and whether there was a tsunami with the magType of mb.

```
quake_mb = quake[quake['magType'] == 'mb'] # to select magType with mb values only

quake_mb_filtered = quake_mb[['mag', 'tsunami']]

correlation_matrix = quake_mb_filtered.corr() # compute correlation

# Plot the heatmap
sns.heatmap(correlation_matrix, annot=True)
plt.title('Correlation: Magnitude vs Tsunami (magType = mb)')
plt.show()
```