

Lecture 06. Plotting and Visualization

Instructor: Luping Yu

Apr 4, 2023

Making informative visualizations (sometimes called plots) is one of the most important tasks in data analysis.

It may be a part of the exploratory process:

- help identify outliers or needed data transformations.
- a way of generating ideas for models.

Python has many add-on libraries for making static or dynamic visualizations, but we will be mainly focused on `matplotlib`. It is a plotting package designed for creating (mostly **two-dimensional**) publication-quality plots.

The project was started in 2002 to enable a *MATLAB-like* plotting interface in Python. `matplotlib` supports various GUI backends on all operating systems and additionally can export visualizations to all of the common vector and raster graphics formats (PDF, SVG, JPG, PNG, BMP, GIF, etc.).

The simplest way to follow the code examples in the chapter is to use interactive plotting in the `Jupyter notebook`. To set this up, execute the following statement in a Jupyter notebook:

```
In [ ]: %matplotlib inline
        %config InlineBackend.figure_format = 'svg'
```

Plotting with pandas

In pandas we may have multiple columns of data, along with row and column labels. `pandas` itself has **built-in** methods that simplify creating visualizations from `DataFrame` and `Series` objects.

Line Plots

`Series` and `DataFrame` each have a `plot` attribute for making some basic plot types. By default, `plot()` `<code> makes line plots:`

```
In [1]: import numpy as np
import pandas as pd

s = pd.Series(np.random.rand(10), index=np.arange(0, 100, 10))
# numpy.random.rand(): 生成随机数
# np.arange(): 生成等差数列

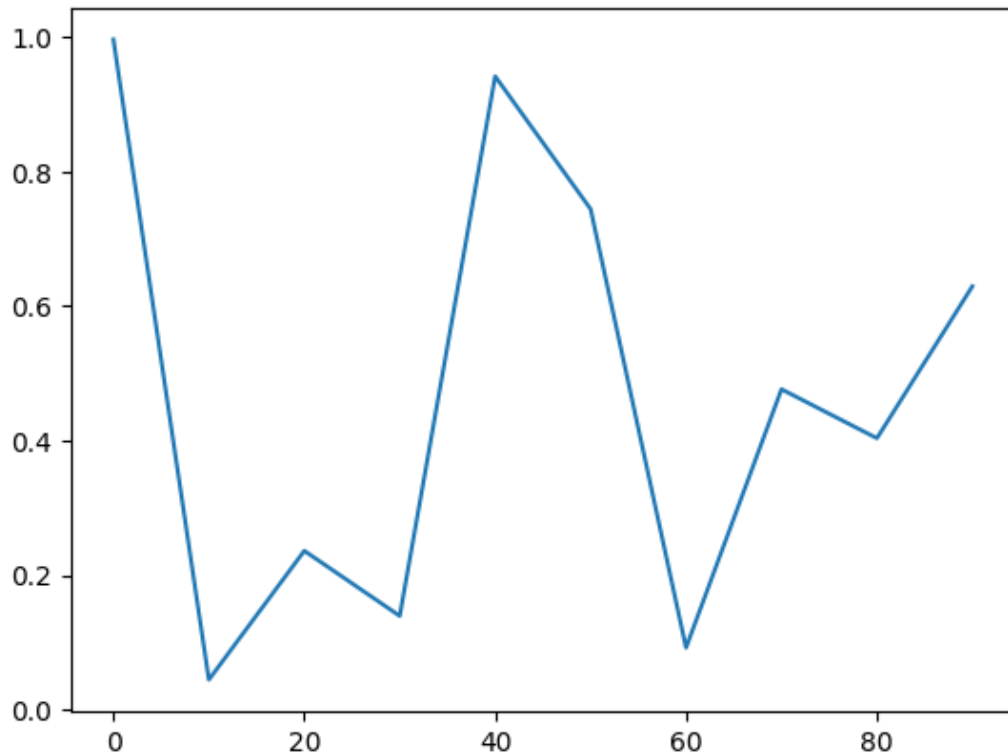
s
```

```
Out[1]: 0      0.996697
      10      0.044160
      20      0.235966
      30      0.139047
      40      0.941851
      50      0.744345
      60      0.091809
      70      0.476299
      80      0.403404
      90      0.629492
      dtype: float64
```

If everything is set up right, a simple line plot should appear:

```
In [2]: s.plot()
```

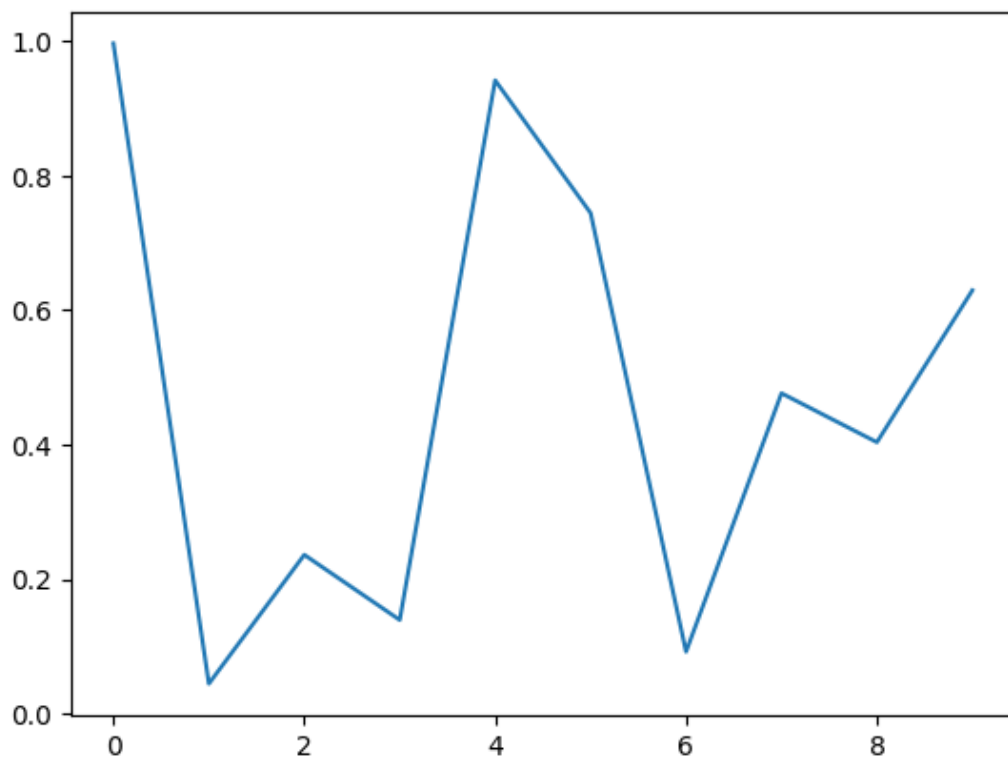
```
Out[2]: <Axes: >
```



The Series object's **index** is passed to matplotlib for plotting on the **x-axis**, though you can disable this by passing `use_index=False`.

```
In [3]: s.plot(use_index=False)
```

```
Out[3]: <Axes: >
```



The x-axis ticks and limits can be adjusted with the `xticks` and `xlim` options, and y-axis respectively with `yticks` and `ylim`. See the following table for a full listing of plot options. I'll comment on a few more of them throughout this section and leave the rest to you to explore.

e.g. `s.plot(xticks=[0,60,100])`, `s.plot(xlim=[60,80])`

- `Series.plot` method arguments:

Argument	Description
<code>ax</code>	matplotlib subplot object to plot on; if nothing passed, uses active matplotlib subplot
<code>style</code>	Style string, like 'ko--', to be passed to matplotlib
<code>alpha</code>	The plot fill opacity (from 0 to 1)
<code>kind</code>	Can be 'area', 'bar', 'barh', 'density', 'hist', 'kde', 'line', 'pie'
<code>logy</code>	Use logarithmic scaling on the y-axis
<code>use_index</code>	Use the object index for tick labels
<code>rot</code>	Rotation of tick labels (0 through 360)
<code>xticks</code>	Values to use for x-axis ticks
<code>yticks</code>	Values to use for y-axis ticks
<code>xlim</code>	x-axis limits (e.g., [0, 10])
<code>ylim</code>	y-axis limits
<code>grid</code>	Display axis grid (on by default)

`DataFrame`'s plot method plots each of its columns as a different line on the same subplot, creating a legend automatically:

```
In [4]: df = pd.DataFrame(np.random.rand(10, 2),
                          columns=['A', 'B'],
                          index=np.arange(0, 100, 10))
```

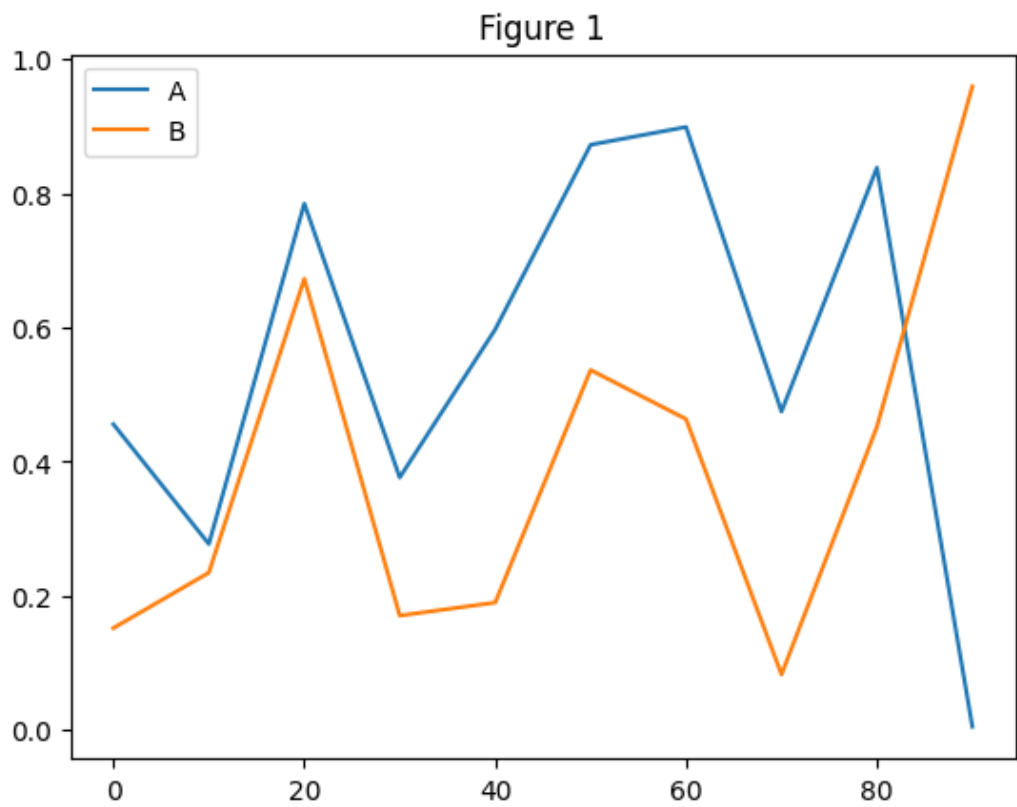
df

Out [4]:

	A	B
0	0.455699	0.151920
10	0.277234	0.234554
20	0.784544	0.672579
30	0.376399	0.170527
40	0.596977	0.189903
50	0.872226	0.536553
60	0.898729	0.463458
70	0.474764	0.082862
80	0.838127	0.451909
90	0.004939	0.959244

```
In [5]: df.plot(title='Figure 1')
```

Out[5]: <Axes: title={'center': 'Figure 1'}>



DataFrame has a number of options allowing some flexibility with how the columns are handled; for example, whether to plot them all on the same subplot or to create separate subplots.

- DataFrame specific plot arguments:

Argument	Description
subplots	Plot each DataFrame column in a separate subplot
sharex	If subplots=True, share the same x-axis, linking ticks and limits
sharey	If subplots=True, share the same y-axis
figsize	Size of figure to create as tuple
title	Plot title as string
legend	Add a subplot legend (True by default)
sort_columns	Plot columns in alphabetical order; by default uses existing column order

Bar Plots

The plot attribute contains a "family" of methods for different plot types. For example, `df.plot()` is equivalent to `df.plot.line()`.

The `plot.bar()` and `plot.barh()` make vertical and horizontal **bar plots**, respectively. In this case, the Series or DataFrame index will be used as the x (bar) or y (barh) ticks:

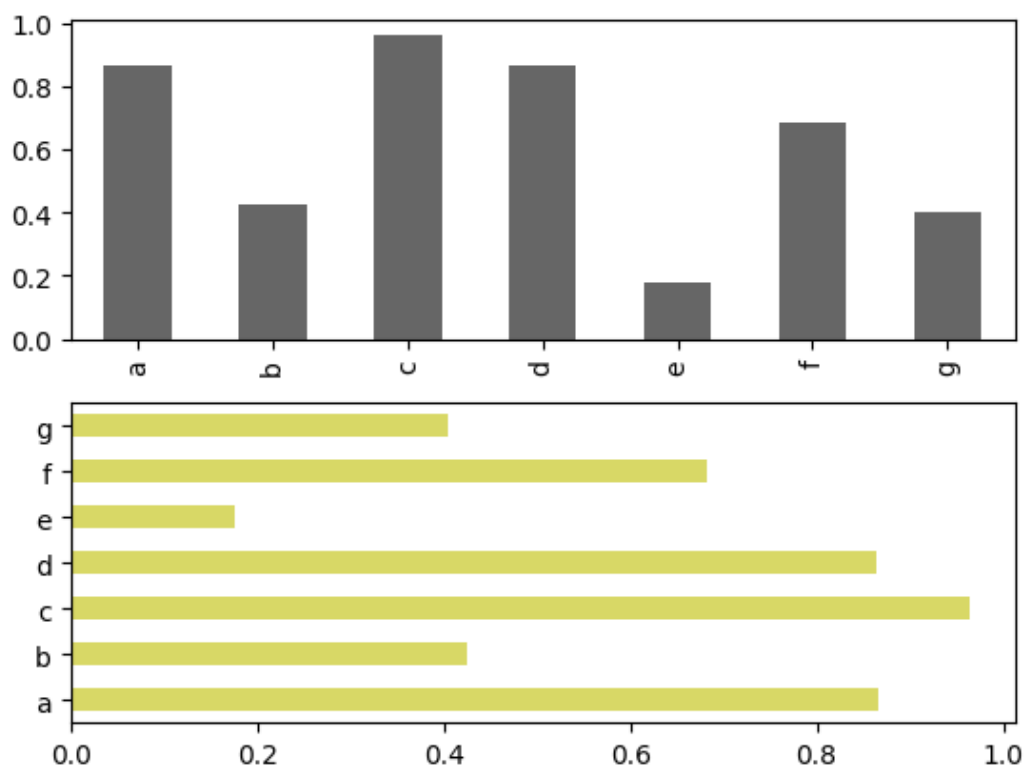
```
In [6]: data = pd.Series(np.random.rand(7), index=list('abcdefg'))
```

data

```
Out[6]: a    0.866845  
b    0.425677  
c    0.964341  
d    0.864273  
e    0.176598  
f    0.682728  
g    0.403872  
dtype: float64
```

```
In [7]: import matplotlib.pyplot as plt  
  
fig, axes = plt.subplots(2, 1) # Subplot with two rows and one column  
  
data.plot.bar(ax=axes[0], color='k', alpha=0.6)  
  
data.plot.barh(ax=axes[1], color='y', alpha=0.6)
```

Out[7]: <Axes: >



The options `color='k'` and `alpha=0.6` set the color of the plots to black and use partial transparency on the filling.

Reference: <https://matplotlib.org/stable/tutorials/colors/colors.html>

With a `DataFrame`, bar plots group the values in each row together in a group in bars, side by side, for each value.

```
In [8]: df = pd.DataFrame(np.random.rand(6, 4),
                          index=['one', 'two', 'three', 'four', 'five', 'six'],
                          columns=pd.Index(['A', 'B', 'C', 'D'], name='Figure 2'))
```

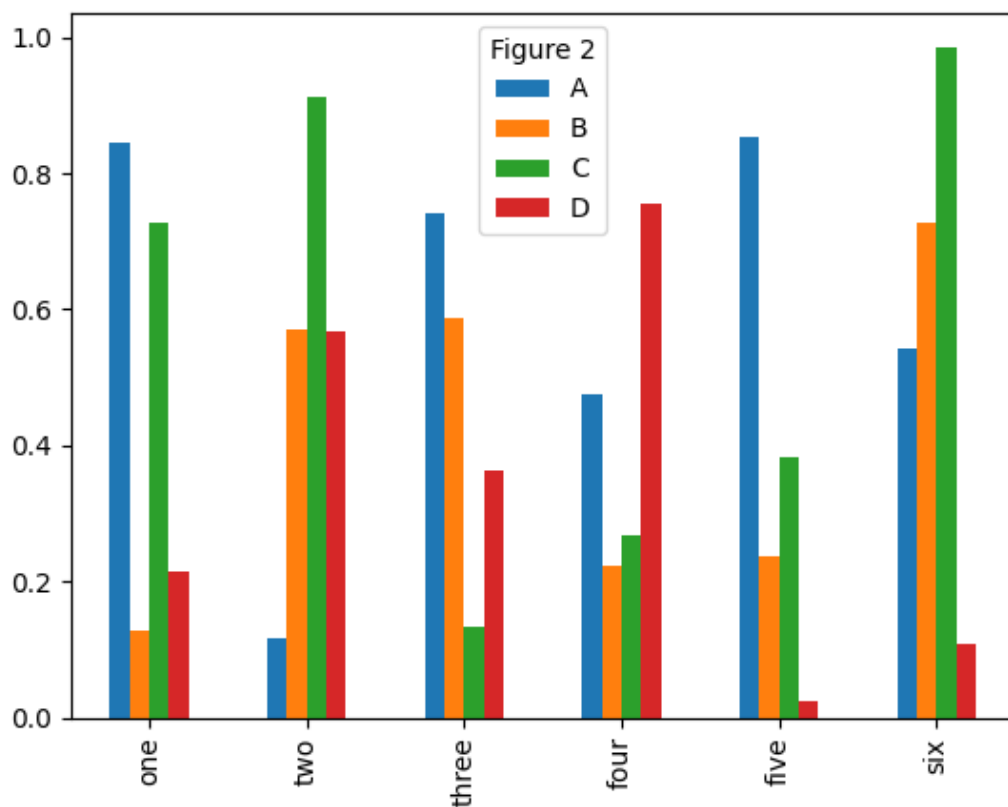
df

```
Out[8]:
```

	Figure 2	A	B	C	D
one		0.844639	0.127464	0.727771	0.215778
two		0.116124	0.571789	0.913120	0.566550
three		0.742316	0.586387	0.132588	0.361745
four		0.476137	0.224144	0.267348	0.755897
five		0.853673	0.235961	0.382031	0.023105
six		0.542335	0.726714	0.986636	0.108582

```
In [9]: df.plot.bar()
```

```
Out[9]: <Axes: >
```

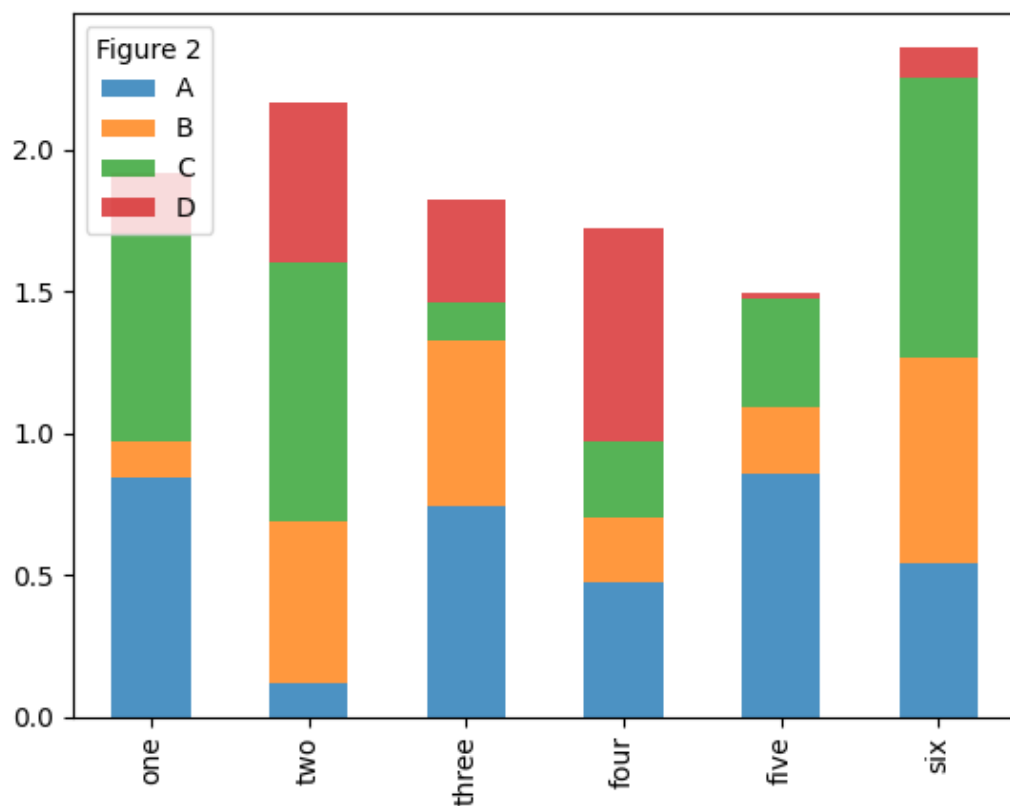


Note that the name "Figure 2" on the DataFrame's columns is used to title the legend.

We create stacked bar plots from a DataFrame by passing `stacked=True`, resulting in the value in each row being stacked together:

```
In [10]: df.plot.bar(stacked=True, alpha=0.8)
```

```
Out[10]: <Axes: >
```



Returning to the `tips.csv` used earlier in [Lecture 04](#), suppose we wanted to make a stacked bar plot showing the percentage of data points for each party size on each day.

We load the data and make a **cross-tabulation** by day and party size:

```
In [11]: tips = pd.read_csv('examples/tips.csv')
```

```
tips
```

```
Out[11]:
```

	total_bill	tip	sex	smoker	day	time	size
0	16.99	1.01	Female	No	Sun	Dinner	2
1	10.34	1.66	Male	No	Sun	Dinner	3
2	21.01	3.50	Male	No	Sun	Dinner	3
3	23.68	3.31	Male	No	Sun	Dinner	2
4	24.59	3.61	Female	No	Sun	Dinner	4
...
239	29.03	5.92	Male	No	Sat	Dinner	3
240	27.18	2.00	Female	Yes	Sat	Dinner	2
241	22.67	2.00	Male	Yes	Sat	Dinner	2
242	17.82	1.75	Male	No	Sat	Dinner	2
243	18.78	3.00	Female	No	Thur	Dinner	2

244 rows × 7 columns

```
In [12]: party_counts = pd.crosstab(tips['day'], tips['size'])
```

```
party_counts
```

```
Out[12]:
```

size	1	2	3	4	5	6
Fri	1	16	1	1	0	0
Sat	2	53	18	13	1	0
Sun	0	39	15	18	3	1
Thur	1	48	4	5	1	3

day

Fri	1	16	1	1	0	0
-----	---	----	---	---	---	---

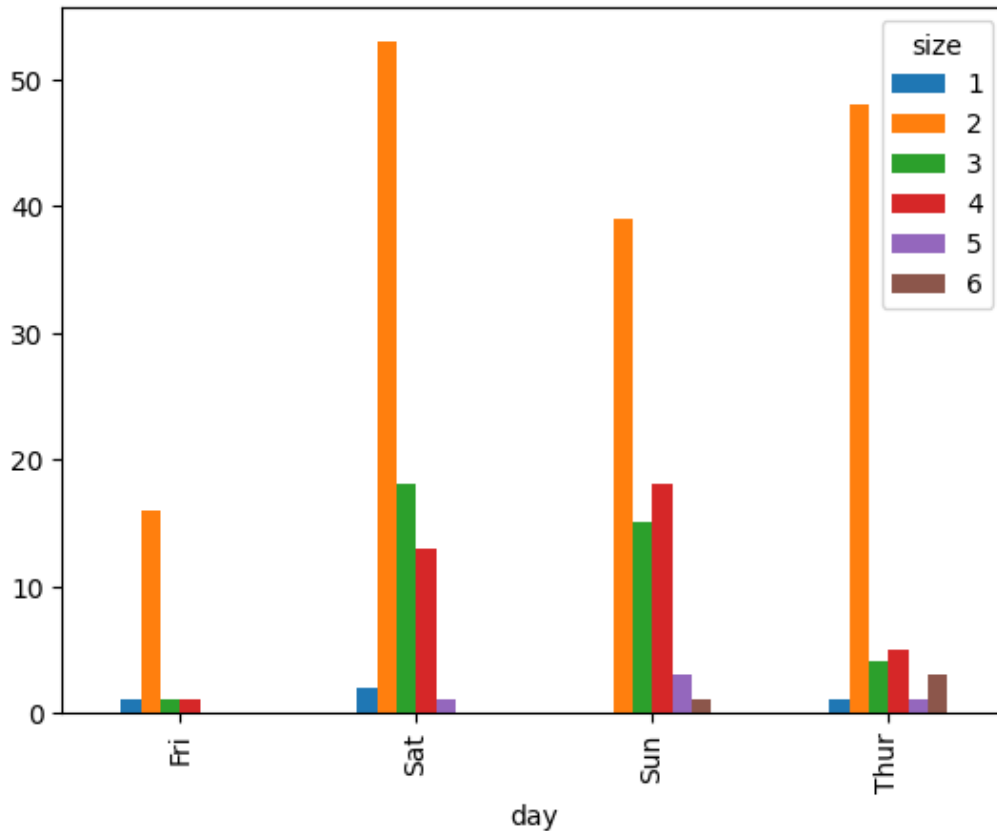
Sat	2	53	18	13	1	0
-----	---	----	----	----	---	---

Sun	0	39	15	18	3	1
-----	---	----	----	----	---	---

Thur	1	48	4	5	1	3
------	---	----	---	---	---	---

```
In [13]: party_counts.plot.bar()
```

```
Out[13]: <Axes: xlabel='day'>
```



So you can see that party sizes appear to increase on the weekend in this dataset.

With data that requires aggregation or summarization before making a plot, using the `seaborn` package can make things much simpler. Let's look now at the tipping percentage by day with seaborn:

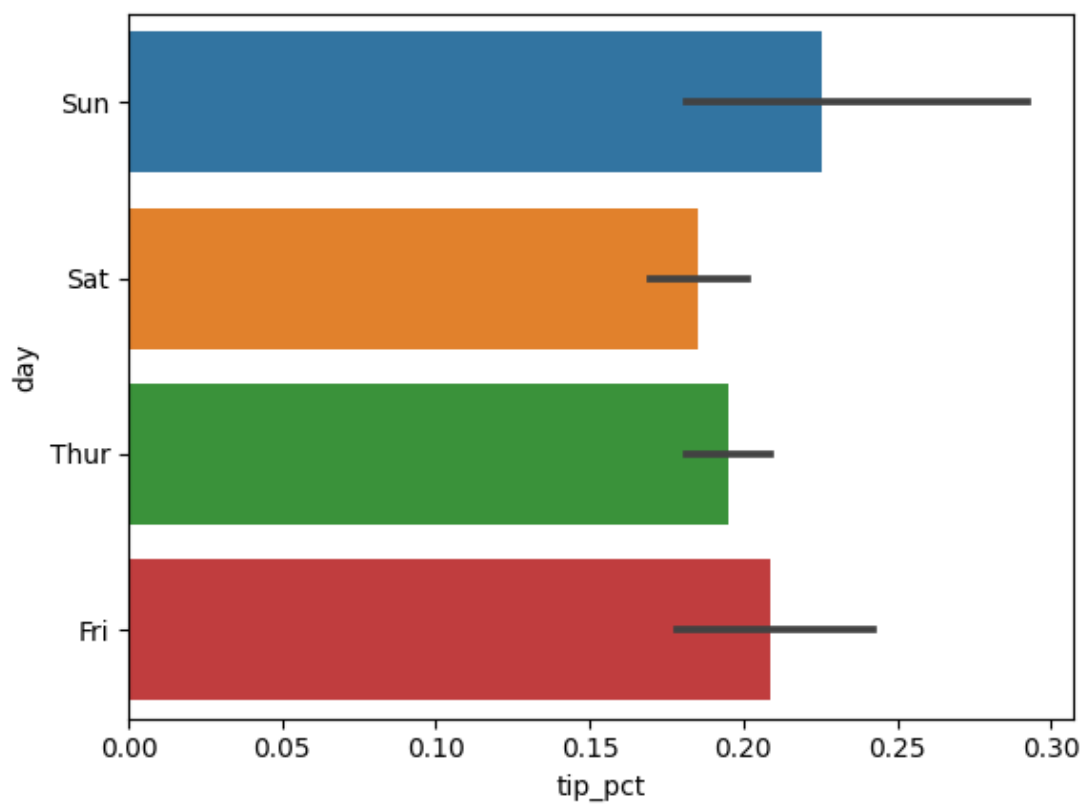
```
In [14]: tips['tip_pct'] = tips['tip'] / (tips['total_bill'] - tips['tip'])
tips.head()
```

```
Out[14]:
```

	total_bill	tip	sex	smoker	day	time	size	tip_pct
0	16.99	1.01	Female	No	Sun	Dinner	2	0.063204
1	10.34	1.66	Male	No	Sun	Dinner	3	0.191244
2	21.01	3.50	Male	No	Sun	Dinner	3	0.199886
3	23.68	3.31	Male	No	Sun	Dinner	2	0.162494
4	24.59	3.61	Female	No	Sun	Dinner	4	0.172069

```
In [15]: import seaborn as sns
sns.barplot(x='tip_pct', y='day', data=tips, orient='h')
```

```
Out[15]: <Axes: xlabel='tip_pct', ylabel='day'>
```

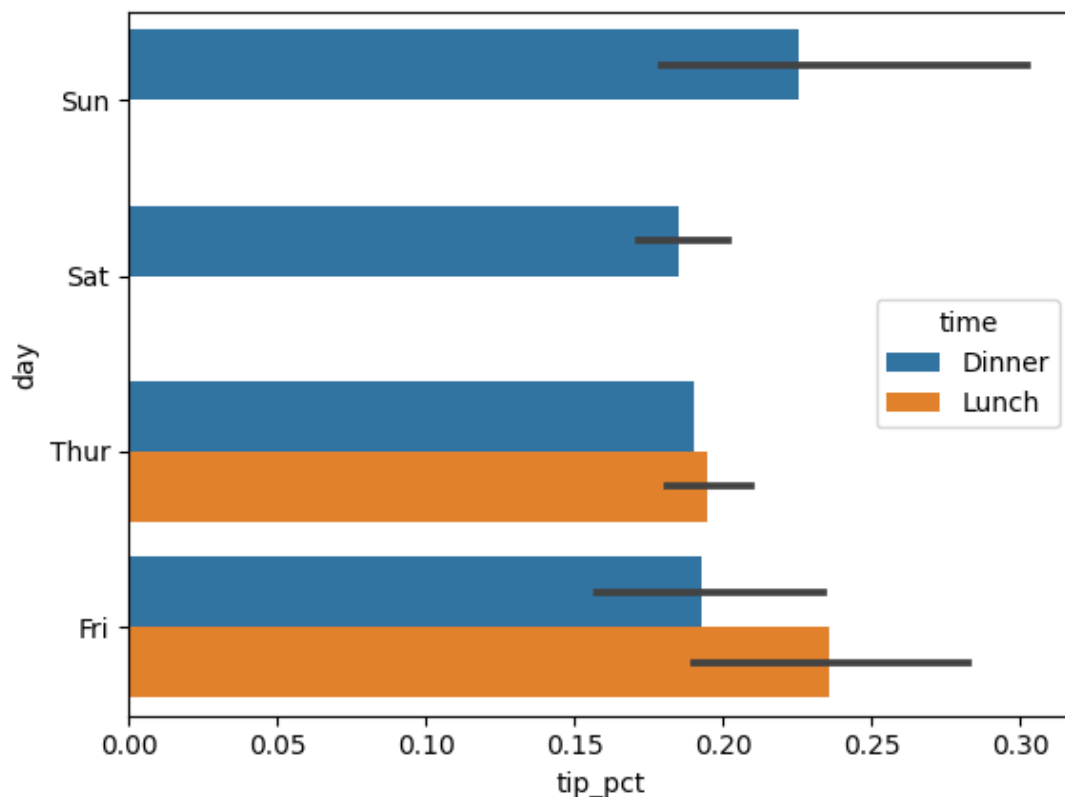
Plotting functions in `seaborn` take a data argument, which can be a pandas DataFrame. The other arguments refer to column names.

Because there are multiple observations for each value in the day, the bars are the **average value** of `tip_pct`. The black lines drawn on the bars represent the **95% confidence interval** (this can be configured through optional arguments).

`seaborn.barplot` has a `hue` option that enables us to split by an additional categorical value:

```
In [16]: sns.barplot(x='tip_pct', y='day', hue='time', data=tips, orient='h')
```

```
Out[16]: <Axes: xlabel='tip_pct', ylabel='day'>
```



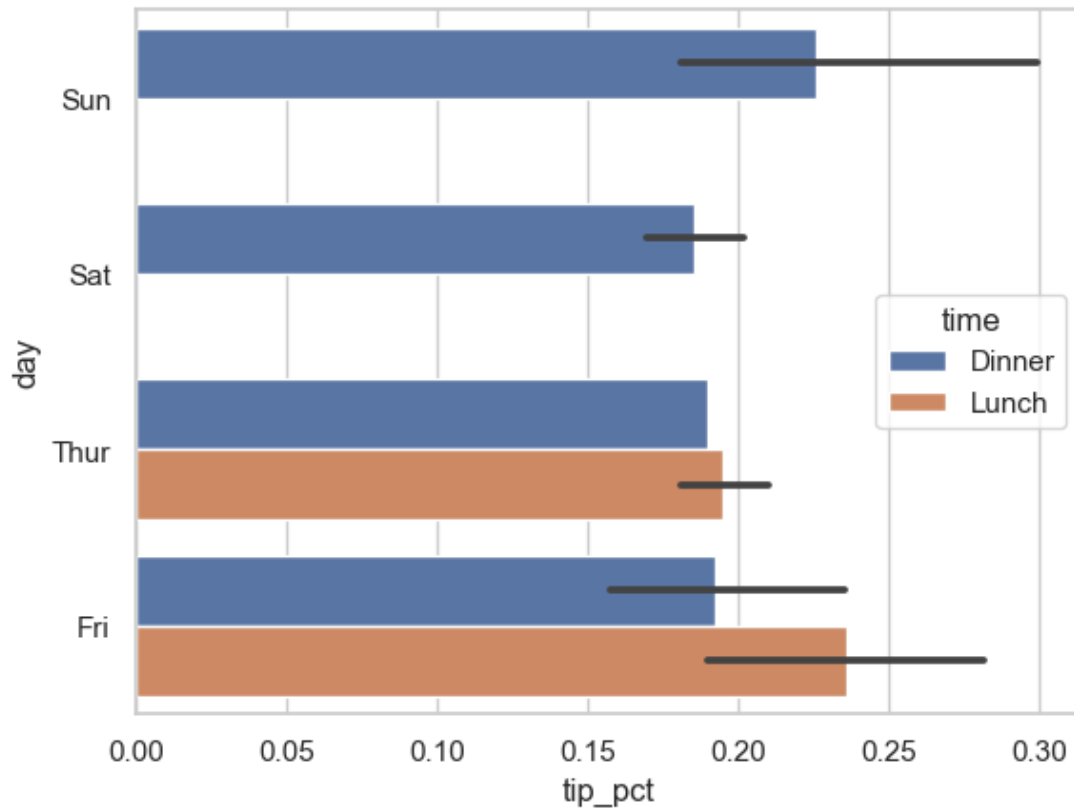
Notice that seaborn has automatically changed the aesthetics of plots: the default color palette, plot background, and grid line colors. You can switch between different plot appearances using

`seaborn.set` :

```
In [17]: sns.set(style='whitegrid')
#sns.reset_orig()

sns.barplot(x='tip_pct', y='day', hue='time', data=tips, orient='h')
```

Out[17]: <Axes: xlabel='tip_pct', ylabel='day'>



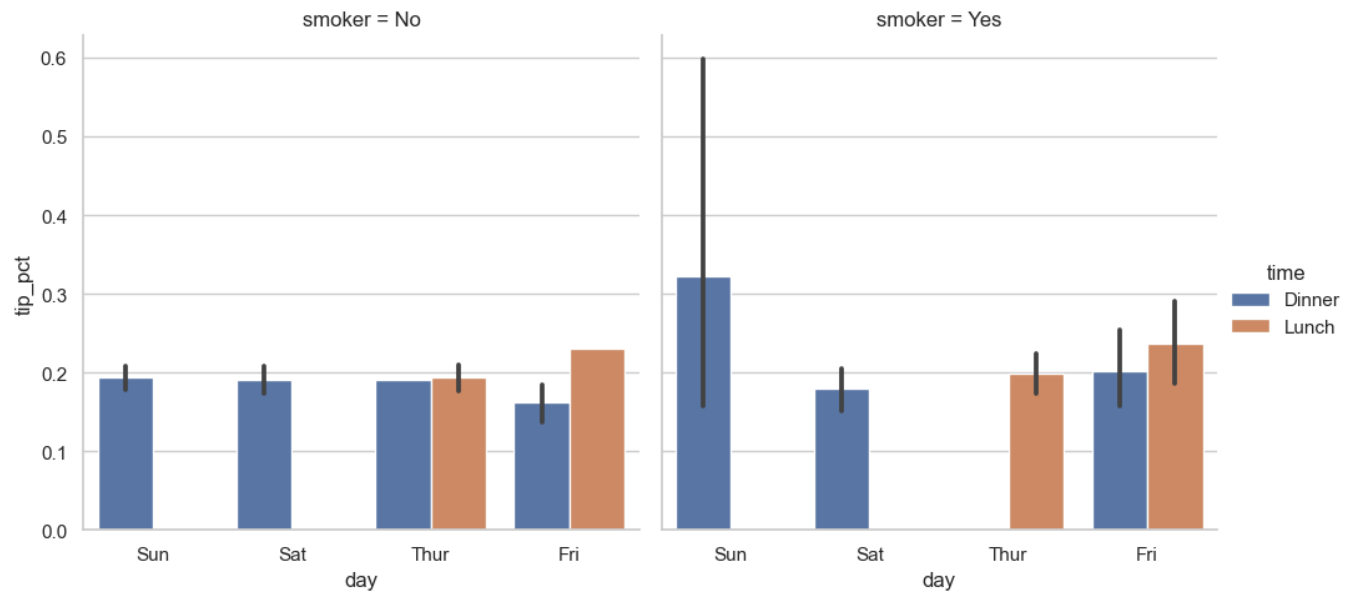
Facet Grids and Categorical Data

What about datasets where we have additional grouping dimensions? One way to visualize data with many categorical variables is to use a `catplot` .

Seaborn has a useful built-in function `catplot` that simplifies making many kinds of plots:

```
In [18]: sns.catplot(x='day', y='tip_pct', hue='time', col='smoker', kind='bar', data=tips)
```

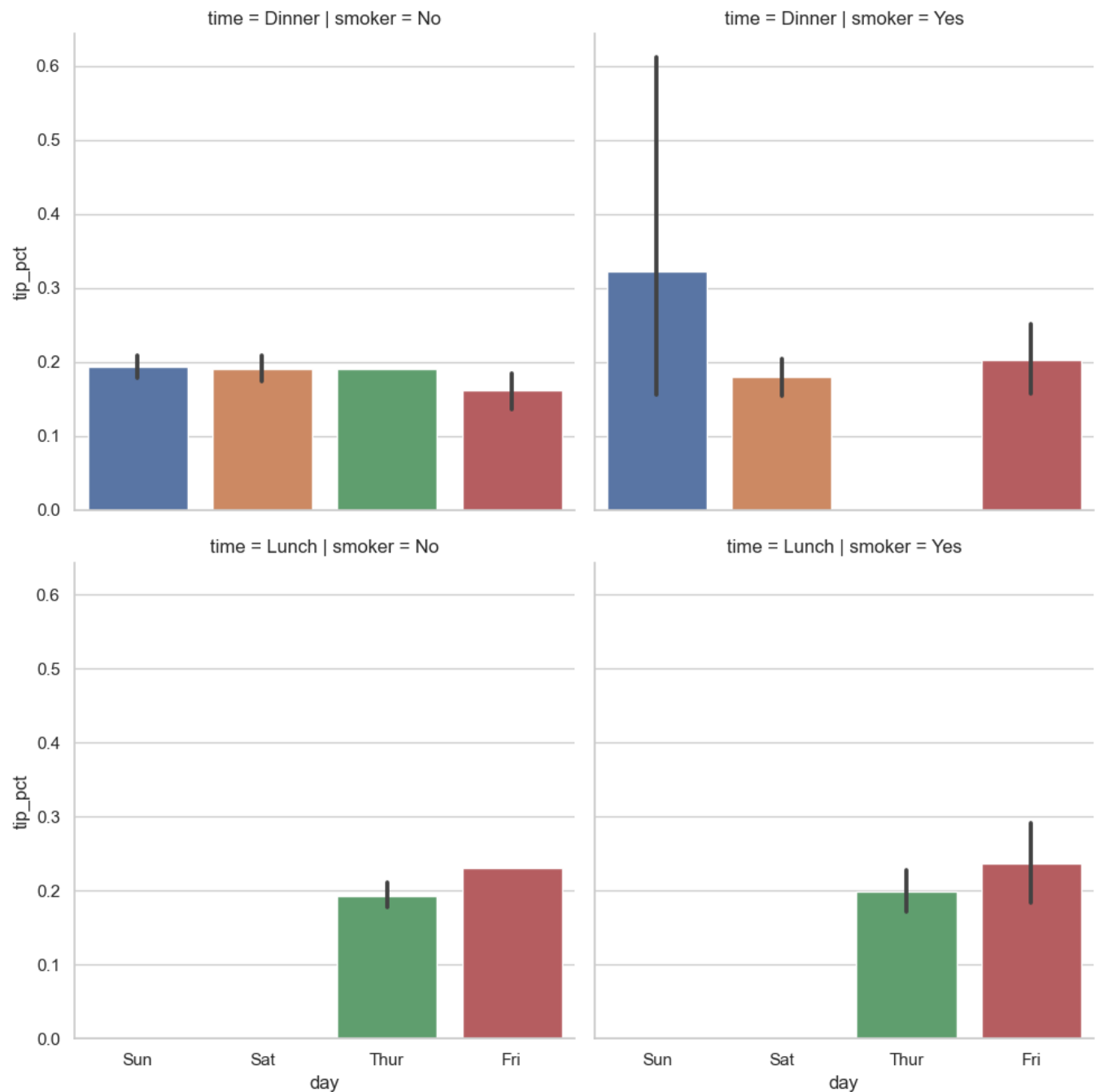
Out[18]: <seaborn.axisgrid.FacetGrid at 0x126f3f510>



Instead of grouping by 'time' by different bar colors within a facet, we can also expand the facet grid by adding one row per time value:

```
In [19]: sns.catplot(x='day', y='tip_pct', row='time', col='smoker', kind='bar', data=tips)
```

```
Out[19]: <seaborn.axisgrid.FacetGrid at 0x126ff0550>
```

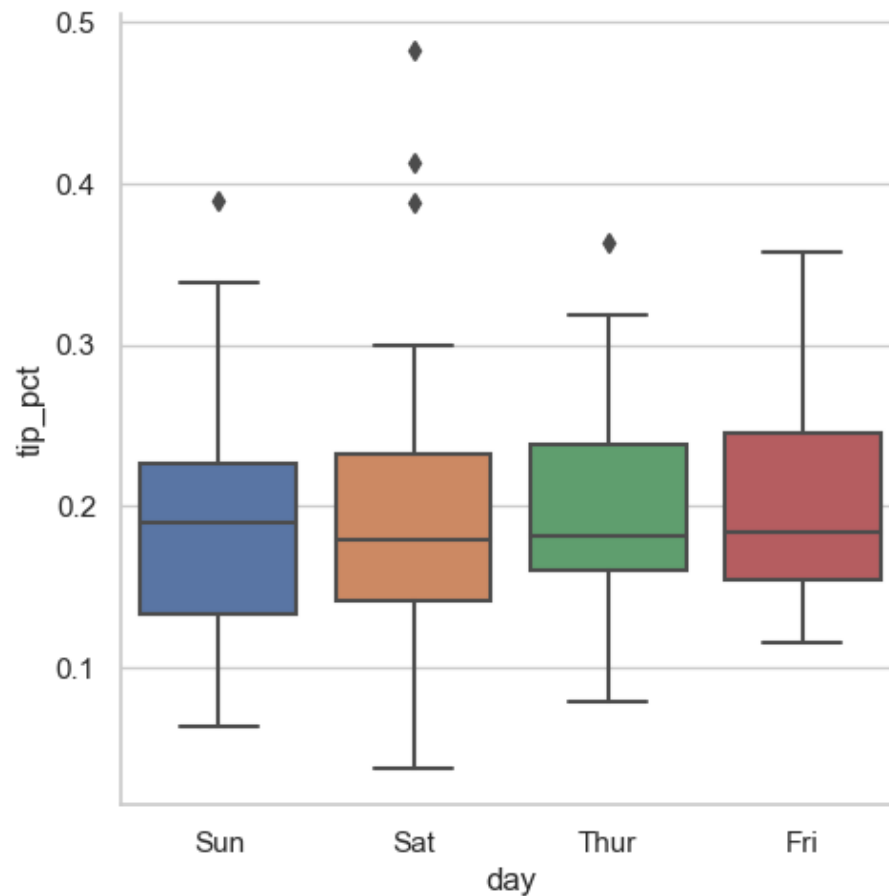


`catplot` supports other plot types that may be useful depending on what you are trying to display.

For example, **box** plots (which show the median, quartiles, and outliers) can be an effective visualization type:

```
In [20]: sns.catplot(x='day', y='tip_pct', kind='box', data=tips[tips.tip_pct < 0.5])
```

```
Out[20]: <seaborn.axisgrid.FacetGrid at 0x126e73750>
```



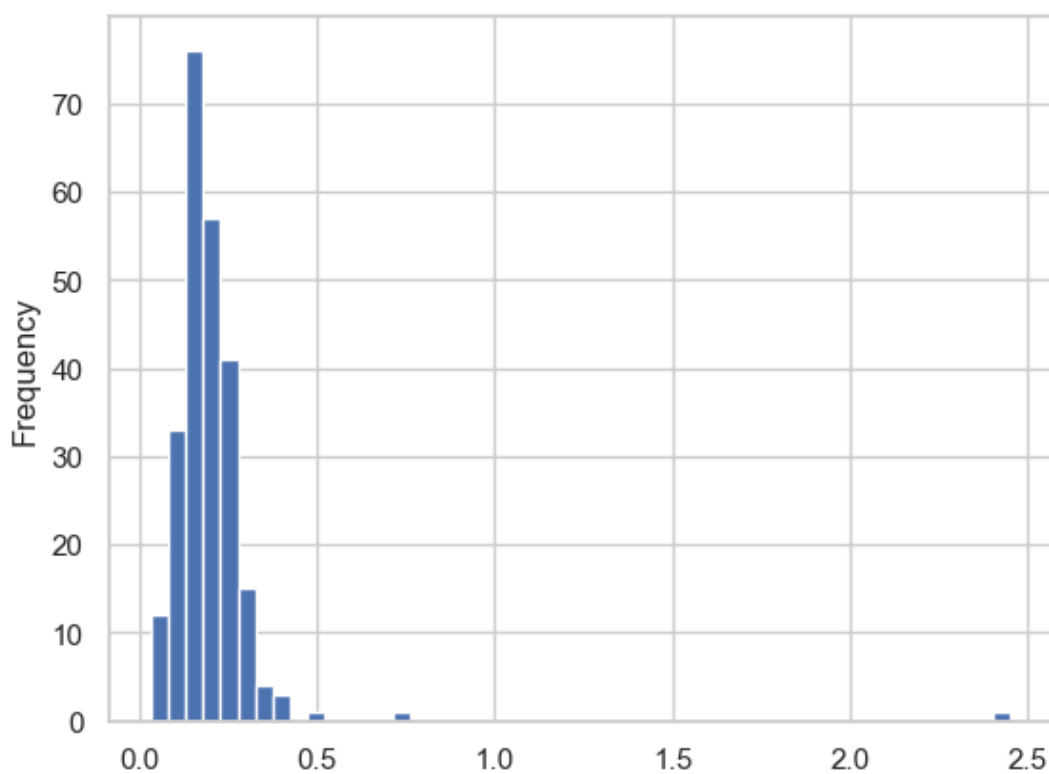
Histograms and Density Plots

A histogram is a kind of bar plot that gives a **discretized display of value frequency**. The data points are split into discrete, evenly spaced bins, and the number of data points in each bin is plotted.

Using the tipping data from before, we can make a histogram of tip percentages of the total bill using the `plot.hist()` method on the Series:

```
In [21]: tips['tip_pct'].plot.hist(bins=50)
```

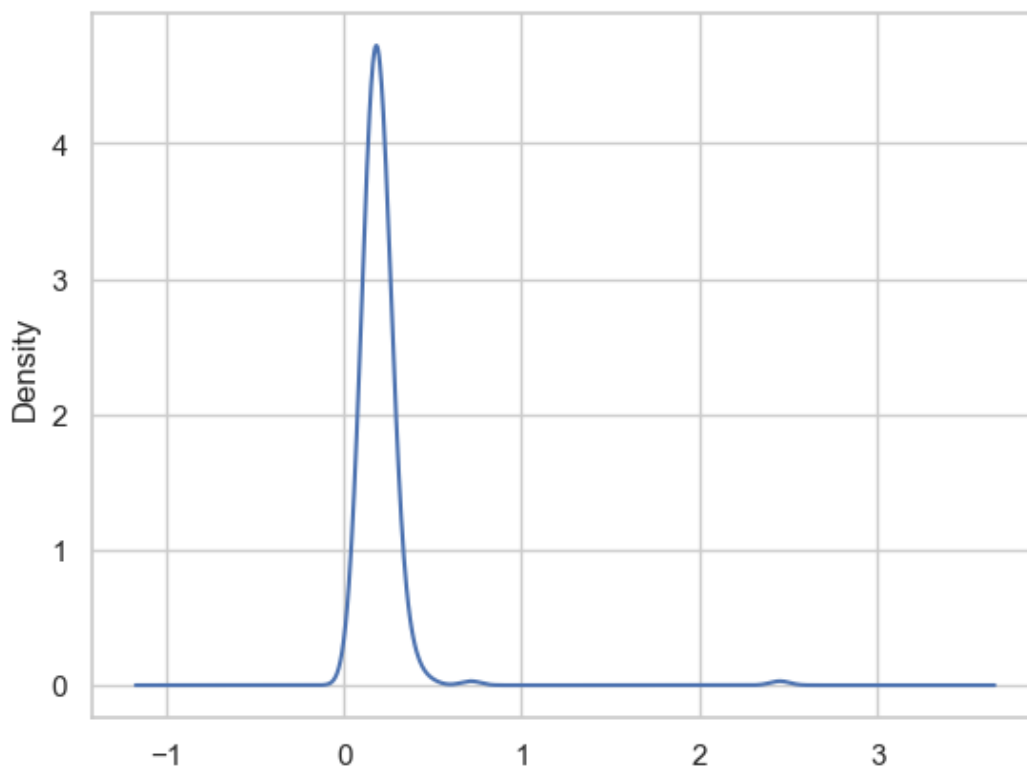
```
Out[21]: <Axes: ylabel='Frequency'>
```



A related plot type is a **density plot**, which is formed by computing an estimate of a continuous probability distribution that might have generated the observed data:

```
In [22]: tips['tip_pct'].plot.density()
```

```
Out[22]: <Axes: ylabel='Density'>
```



Seaborn makes histograms and density plots even easier through its `.histplot()` method, which can plot both a histogram and a continuous density estimate simultaneously.

As an example, consider a bimodal distribution consisting of draws from two different standard normal distributions:

```
In [23]: comp1 = np.random.normal(0, 1, size=200)
         comp2 = np.random.normal(10, 2, size=200)
```

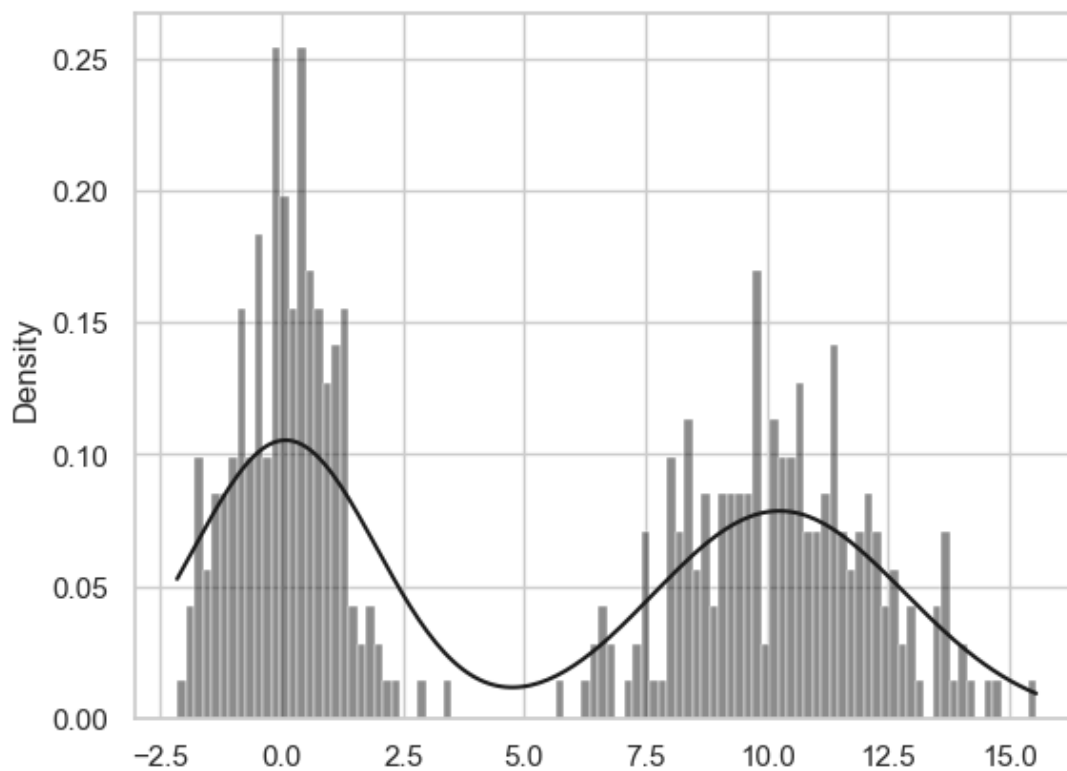
```
values = pd.Series(np.concatenate([comp1, comp2]))
```

```
values
```

```
Out[23]: 0      -0.161691
         1       0.359412
         2       0.802309
         3       0.952121
         4      -0.013733
         ...
        395     10.375722
        396       8.726607
        397       9.719123
        398     12.276264
        399     11.101543
        Length: 400, dtype: float64
```

```
In [24]: sns.histplot(values, bins=100, color='k', kde=True, stat="density")
```

```
Out[24]: <Axes: ylabel='Density'>
```



Scatter or Point Plots

Point plots or scatter plots can be a useful way of examining the relationship between two one-dimensional data series.

For example, here we load the [macrodata.csv](#), select a few variables, then compute log differences:

```
In [25]: macro = pd.read_csv('examples/macrodata.csv')

macro
```

Out[25]:

	year	quarter	realgdp	realcons	realinv	realgovt	realdpi	cpi	m1	tbilrate	unemp	
0	1959.0	1.0	2710.349	1707.4	286.898	470.045	1886.9	28.980	139.7	2.82	5.8	17
1	1959.0	2.0	2778.801	1733.7	310.859	481.301	1919.7	29.150	141.7	3.08	5.1	17
2	1959.0	3.0	2775.488	1751.8	289.226	491.260	1916.4	29.350	140.5	3.82	5.3	17
3	1959.0	4.0	2785.204	1753.7	299.356	484.052	1931.3	29.370	140.0	4.33	5.6	17
4	1960.0	1.0	2847.699	1770.5	331.722	462.199	1955.5	29.540	139.6	3.50	5.2	18
...
198	2008.0	3.0	13324.600	9267.7	1990.693	991.551	9838.3	216.889	1474.7	1.17	6.0	30
199	2008.0	4.0	13141.920	9195.3	1857.661	1007.273	9920.4	212.174	1576.5	0.12	6.9	30
200	2009.0	1.0	12925.410	9209.2	1558.494	996.287	9926.4	212.671	1592.8	0.22	8.1	30
201	2009.0	2.0	12901.504	9189.0	1456.678	1023.528	10077.5	214.469	1653.6	0.18	9.2	30
202	2009.0	3.0	12990.341	9256.0	1486.398	1044.088	10040.6	216.385	1673.9	0.12	9.6	30

203 rows × 14 columns

In [26]:

```
data = macro[['cpi', 'm1', 'tbilrate', 'unemp']]

trans_data = np.log(data).diff().dropna()
# np.log(): Natural logarithm
# diff(): First discrete difference of element

trans_data
```

Out[26]:

	cpi	m1	tbilrate	unemp
1	0.005849	0.014215	0.088193	-0.128617
2	0.006838	-0.008505	0.215321	0.038466
3	0.000681	-0.003565	0.125317	0.055060
4	0.005772	-0.002861	-0.212805	-0.074108
5	0.000338	0.004289	-0.266946	0.000000
...
198	-0.007904	0.045361	-0.396881	0.105361
199	-0.021979	0.066753	-2.277267	0.139762
200	0.002340	0.010286	0.606136	0.160343
201	0.008419	0.037461	-0.200671	0.127339
202	0.008894	0.012202	-0.405465	0.042560

202 rows × 4 columns

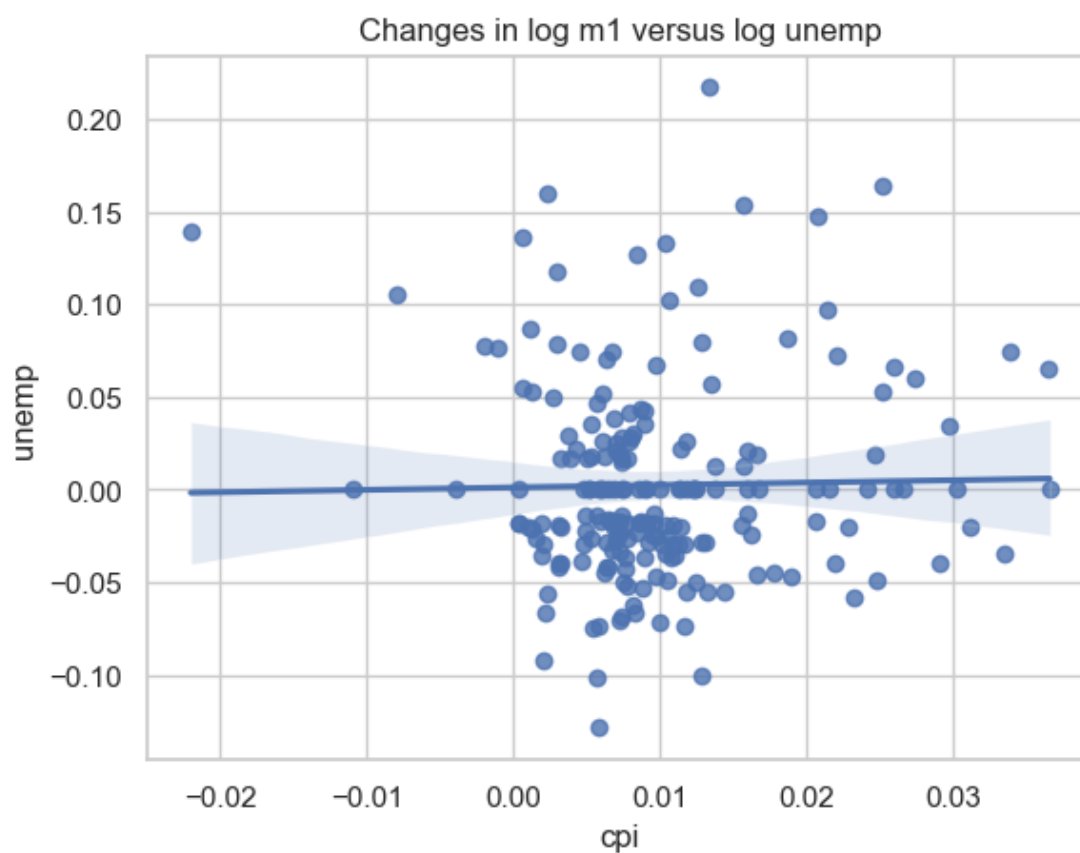
We can then use seaborn's `regplot()` method, which makes a scatter plot and fits a linear regression line:

In [27]:

```
sns.regplot(x='cpi', y='unemp', data=trans_data)

plt.title('Changes in log %s versus log %s' % ('m1', 'unemp'))
```

Out[27]: Text(0.5, 1.0, 'Changes in log m1 versus log unemp')



Conclusion

The goal of this chapter was to get your feet wet with some basic data visualization using **pandas**, **matplotlib**, and **seaborn**.

If visually communicating the results of data analysis is important in your work, I encourage you to seek out resources to learn more about effective data visualization.

It is an active field of research and you can practice with many excellent learning resources available online and in print form.

Example: plot of the World Population

```
In [28]: import pandas as pd
import plotly.offline as offline
```

```
In [29]: df = pd.read_csv('examples/worldbank_population.csv', encoding = 'ISO-8859-1')

# The sunburst plot requires weights (values), labels, and parent (region, or World)
# We build the corresponding table here
columns = ['parents', 'labels', 'values']

level1 = df.copy()
level1.columns = columns
level1['text'] = level1['values'].apply(lambda pop: '{:,.0f}'.format(pop))

level2 = df.groupby('region').population.sum().reset_index()[['region', 'region', 'population']]
level2.columns = columns
level2['parents'] = 'World'
# move value to text for this level
level2['text'] = level2['values'].apply(lambda pop: '{:,.0f}'.format(pop))
level2['values'] = 0

level3 = pd.DataFrame({'parents': [''], 'labels': ['World'],
                       'values': [0.0], 'text': ['{:,.0f}'.format(df['population'].sum())]})
```



```
all_levels = pd.concat([level1, level2, level3], axis=0).reset_index(drop=True)
all_levels
```

Out[29]:

	parents	labels	values	text
0	Latin America & Caribbean	Aruba	106537.0	106,537
1	South Asia	Afghanistan	40099462.0	40,099,462
2	Sub-Saharan Africa	Angola	34503774.0	34,503,774
3	Europe & Central Asia	Albania	2811666.0	2,811,666
4	Europe & Central Asia	Andorra	79034.0	79,034
...
221	World	Middle East & North Africa	0.0	486,167,363
222	World	North America	0.0	370,203,720
223	World	South Asia	0.0	1,901,911,604
224	World	Sub-Saharan Africa	0.0	1,181,162,739
225		World	0.0	7,864,921,177

226 rows × 4 columns

In [30]:

```
# And now we can plot the World Population
offline.ipplot(dict(
    data=[dict(type='sunburst', hoverinfo='text', **all_levels)],
    layout=dict(title='World Population (World Bank, 2023)<br>Click on a region to zoom',
        width=800, height=800)),
    validate=False)
```

World Population (World Bank, 2023)
Click on a region to zoom

