

Table of Contents

Introduction	1.1
1. The Data Science Lifecycle	1.2
1.1 About This Book	1.2.1
1.2 The Students of Data 100	1.2.2
1.3 Exploring the Data	1.2.3
1.4 What's in a Name?	1.2.4
2. Data Generation	1.3
2.1 Dewey Defeats Truman	1.3.1
2.2 Probability Sampling	1.3.2
2.3 SRS vs. "Big Data"	1.3.3
3. Tabular Data and pandas	1.4
3.1 Indexes, Slicing, and Sorting	1.4.1
3.2 Grouping and Pivoting	1.4.2
3.3 Apply, Strings, and Plotting	1.4.3
4. Data Cleaning	1.5
4.1 Cleaning the Calls Dataset	1.5.1
4.2 Cleaning the Stops Dataset	1.5.2
5. Exploratory Data Analysis	1.6
5.1 Structure and Joins	1.6.1
5.2 Granularity	1.6.2
5.3 Scope	1.6.3
5.4 Temporality	1.6.4
5.5 Faithfulness	1.6.5
6. Data Visualization	1.7
6.1 Quantitative Data	1.7.1
6.2 Qualitative Data	1.7.2
6.3 Customizing Plots	1.7.3
6.4 Principles of Visualization	1.7.4
6.5 Principles of Visualization 2	1.7.5
6.6 Visualization Philosophy	1.7.6

7. Web Technologies [In Progress]	1.8
8. Working With Text [In Progress]	1.9
9. Databases and SQL [In Progress]	1.10
10. Modeling and Estimation	1.11
10.1 A Simple Model	1.11.1
10.2 Cost Functions	1.11.2
10.3 Absolute Cost and Huber Cost	1.11.3
11. Gradient Descent	1.12
11.1 Basic Numerical Optimization	1.12.1
11.2 Defining Gradient Descent	1.12.2
12. Linear Regression	1.13
12.1 Defining a Simple Linear Model	1.13.1
12.2 Fitting the Model	1.13.2
12.3 Multiple Linear Regression	1.13.3

Principles and Techniques of Data Science

By [Sam Lau](#), [Joey Gonzalez](#), and [Deb Nolan](#).

This is the textbook for [Data 100](#), the Principles and Techniques of Data Science course at UC Berkeley.

DS100 is the upper-division data science course intended to be taken after [Data 8](#), the Foundations of Data Science course.

The contents of this book are licensed for free consumption under the following license:
[Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International \(CC BY-NC-ND 4.0\)](#)

The Data Science Lifecycle

In data science, we use large and diverse data sets to make conclusions about the world. In this course, we will discuss the key principles and techniques of data science through the dual lens of computational and inferential thinking. Practically speaking, this involves the following process:

1. Formulating a question or problem
2. Acquiring and cleaning data
3. Conducting exploratory data analysis
4. Using prediction and inference to draw conclusions

It is quite common for more questions and problems to emerge after the last step of this process, and we can thus repeatedly engage in this procedure to discover new characteristics of our world. This positive feedback loop is so central to our work that we call it the **data science lifecycle**.

If the data science lifecycle were as easy to conduct as it is to state, there would be no need for textbooks about the subject. Fortunately, each of the steps in the lifecycle contain numerous challenges that reveal powerful and often surprising insights that form the foundation of making thoughtful decisions using data.

As in Data 8, we will begin with an example.

About This Book

Before we continue, it is important for us to state our assumptions about the reader.

In this book, we will proceed as though you have taken [Data 8](#) or some equivalent. In particular, we will assume that you have some familiarity with the following topics (links to pages from the Data 8 textbook are given in parentheses).

- Tabular data manipulation: selection, filtering, grouping, joining ([link](#))
- Sampling, empirical distributions of statistics ([link](#))
- Hypothesis testing using bootstrap resampling ([link](#))
- Least squares regression and regression inference ([link](#))
- Classification ([link](#))

In addition, we will assume that you have taken [CS61A](#) or some equivalent and thus will not explain Python syntax except in special cases.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [The Students of Data 100](#)
- [Question Formulation](#)
- [Data Acquisition and Cleaning](#)

The Students of Data 100

Recall that the data science lifecycle involves the following broad steps:

1. Question/Problem Formulation:

- i. What do we want to know or what problems are we trying to solve?
- ii. What are our hypotheses?
- iii. What are our metrics of success?

2. Data Acquisition and Cleaning:

- i. What data do we have and what data do we need?
- ii. How will we collect more data?
- iii. How do we organize the data for analysis?

3. Exploratory Data Analysis:

- i. Do we already have relevant data?
- ii. What are the biases, anomalies, or other issues with the data?
- iii. How do we transform the data to enable effective analysis?

4. Prediction and Inference:

- i. What does the data say about the world?
- ii. Does it answer our questions or accurately solve the problem?
- iii. How robust are our conclusions?

Question Formulation

We would like to figure out if the data we have on student names in Data 100 give us any additional information about the students themselves. Although this is a vague question to ask, it is enough to get us working with our data and we can surely make the question more

precise as we go.

Data Acquisition and Cleaning

In Data 100, we will study various methods to collect data.

Let's begin by looking at our data, the roster of student first names that we've downloaded from a previous offering of Data 100.

Don't worry if you don't understand the code for now; we'll introduce the libraries in more depth later. Instead, focus on the process and the charts that we show.

```
import pandas as pd

students = pd.read_csv('roster.csv')
students
```

	Name	Role
0	Keeley	Student
1	John	Student
2	BRYAN	Student
...
276	Ernesto	Waitlist Student
277	Athan	Waitlist Student
278	Michael	Waitlist Student

279 rows × 2 columns

We can quickly see that there are some quirks in the data. For example, one of the student's names is all uppercase letters. In addition, it is not obvious what the Role column is for.

In Data 100, we will study how to identify anomalies in data and apply corrections. The differences in capitalization will cause our programs to think that 'BRYAN' and 'Bryan' are different names when they are identical for our purposes. Let's convert all names to lower case to avoid this.

```
students['Name'] = students['Name'].str.lower()
students
```

1.2 The Students of Data 100

	Name	Role
0	keeley	Student
1	john	Student
2	bryan	Student
...
276	ernesto	Waitlist Student
277	athan	Waitlist Student
278	michael	Waitlist Student

279 rows × 2 columns

Now that our data are in a format that's easier for us to work with, let's proceed to exploratory data analysis.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Exploratory Data Analysis](#)
- [What's in a Name?](#)

Exploratory Data Analysis

The term Exploratory Data Analysis (EDA for short) refers to the process of discovering traits about our data that inform future analysis.

Here's the `students` table from the previous page:

students		
	Name	Role
0	keeley	Student
1	john	Student
2	bryan	Student
...
276	ernesto	Waitlist Student
277	athan	Waitlist Student
278	michael	Waitlist Student

279 rows × 2 columns

We are left with a number of questions. How many students are in this roster? What does the `Role` column mean? We conduct EDA in order to understand our data more thoroughly.

In Data 100 we will study exploratory data analysis and practice analyzing new datasets.

Oftentimes, we explore the data by repeatedly posing simple questions about the data that we'd like to know about. We will structure our analysis this way.

How many students are in our dataset?

```
print("There are", len(students), "students on the roster.")
```

There are 279 students on the roster.

A natural follow-up question is whether this is the complete list of students or not. In this case, we happen to know that this list contains all students in the class.

What is the meaning of the `Role` field?

Understanding the meaning of field can often be achieved by looking at the unique values of the field's data:

```
students['Role'].value_counts().to_frame()
```

	Role
Student	237
Waitlist Student	42

We can see here that our data contain not only students enrolled in the class at the time but also the students on the waitlist. The `Role` column tells us whether each student is enrolled.

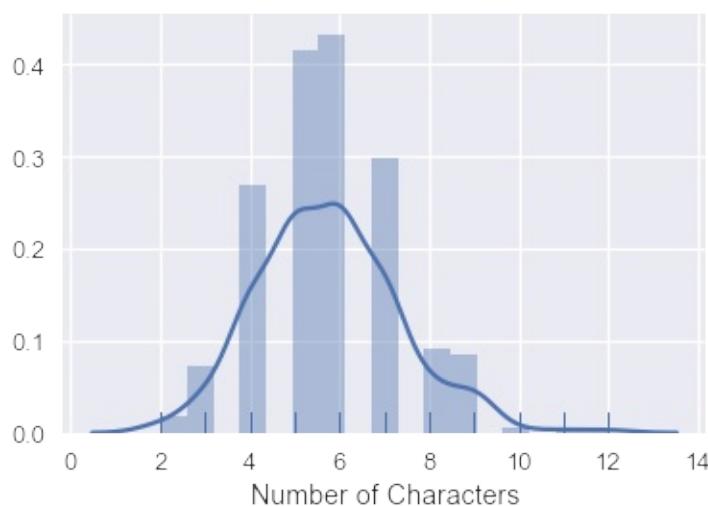
What about the names? How can we summarize this field?

In Data 100 we will deal with many different kinds of data (not just numbers) and we will study techniques to diverse types of data.

A good starting point might be to examine the lengths of the strings.

```
sns.distplot(students['Name'].str.len(), rug=True,
             xlabel="Number of Characters")
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x10e6fd0b8>
```



This visualization shows us that most names are between 3 and 9 characters long. This gives us a chance to check whether our data seem reasonable — if there were many names that were 1 character long we'd have good reason to re-examine our data.

What's in a Name?

Although this dataset is rather simple, we will soon see that first names alone can reveal quite a bit about our class.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [What's in a Name?](#)
 - [Inferring Sex From Name](#)
 - [Inferring Age From Name](#)

What's in a Name?¶

So far, we have asked a broad question about our data: "Do the first names of students in Data 100 tell us anything about the class?"

We have done a bit of data cleaning by converting all our names to lowercase. During our exploratory data analysis we discovered that our roster contains about 270 names of students in the class and on the waitlist, and that most of our first names are between 4 and 8 characters long.

What else can we discover about our class based on first names? We might consider a single name from our dataset:

```
students['Name'][5]
```

```
'jerry'
```

From this name we can infer that the student is likely a male. We can also take a guess at the student's age. For example, if we happen to know that Jerry was a very popular baby name in 1998, we might guess that this student is around twenty years old.

This thinking gives us two new questions to investigate:

1. "Do the first names of students in Data 100 tell us the distribution of sex in the class?"
2. "Do the first names of students in Data 100 tell us the distribution of ages in the class?"

In order to investigate these questions, we will need a dataset that associates names with sex and year. Conveniently, the US Social Security department hosts such a dataset online: <https://www.ssa.gov/oact/babynames/index.html>. Their dataset records the names given to babies at birth and as such is often referred to as the Baby Names dataset.

We will start by downloading and then loading the dataset into Python. Again, don't worry about understanding the code in this first chapter. It's more important that you understand the overall process.

```
import urllib.request
import os.path

data_url = "https://www.ssa.gov/oact/babynames/names.zip"
local_filename = "babynames.zip"
if not os.path.exists(local_filename): # if the data exists
    don't download again
    with urllib.request.urlopen(data_url) as resp,
        open(local_filename, 'wb') as f:
            f.write(resp.read())

import zipfile
babynames = []
with zipfile.ZipFile(local_filename, "r") as zf:
    data_files = [f for f in zf.filelist if f.filename[-3:] ==
        "txt"]
    def extract_year_from_filename(fn):
        return int(fn[3:7])
    for f in data_files:
        year = extract_year_from_filename(f.filename)
        with zf.open(f) as fp:
            df = pd.read_csv(fp, names=["Name", "Sex", "Count"])
            df["Year"] = year
            babynames.append(df)
babynames = pd.concat(babynames)
babynames
```

1.4 What's in a Name?

	Name	Sex	Count	Year
0	Mary	F	9217	1884
1	Anna	F	3860	1884
2	Emma	F	2587	1884
...
2081	Verna	M	5	1883
2082	Winnie	M	5	1883
2083	Winthrop	M	5	1883

1891894 rows × 4 columns

```
ls -alh babynames.csv
```

```
-rw-r--r-- 1 sam staff 30M Jan 22 15:31 babynames.csv
```

It looks like the dataset contains names, the sex given to the baby, the number of babies with that name, and the year of birth for those babies. To be sure, we check the dataset description from the SSN Office: <https://www.ssa.gov/oact/babynames/background.html>

All names are from Social Security card applications for births that occurred in the United States after 1879. Note that many people born before 1937 never applied for a Social Security card, so their names are not included in our data. For others who did apply, our records may not show the place of birth, and again their names are not included in our data.

All data are from a 100% sample of our records on Social Security card applications as of March 2017.

A useful visualization for this data is to plot the number of male and female babies born each year:

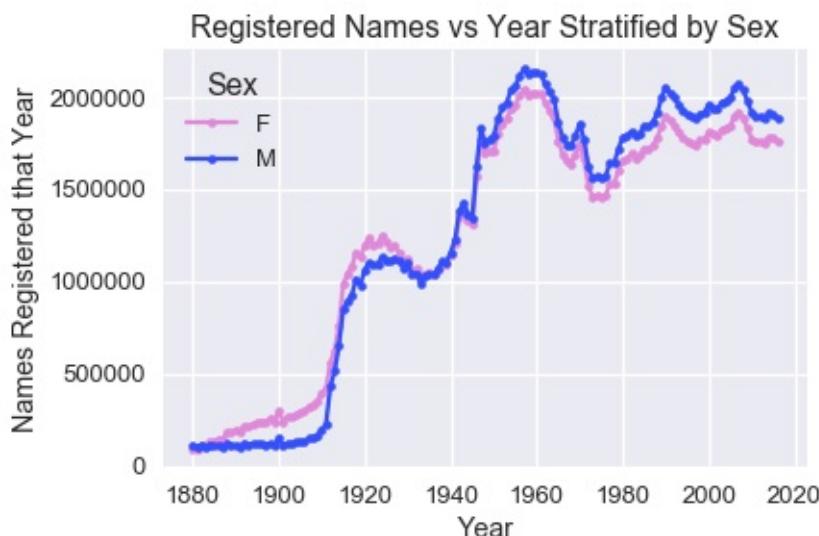
```

pivot_year_name_count = pd.pivot_table(
    babynames, index='Year', columns='Sex',
    values='Count', aggfunc=np.sum)

pink_blue = ["#E188DB", "#334FFF"]
with sns.color_palette(sns.color_palette(pink_blue)):
    pivot_year_name_count.plot(marker=".")  

    plt.title("Registered Names vs Year Stratified by Sex")
    plt.ylabel('Names Registered that Year')

```



This plot makes us question whether the US had babies in 1880. A sentence from the quote above helps explain:

Note that many people born before 1937 never applied for a Social Security card, so their names are not included in our data. For others who did apply, our records may not show the place of birth, and again their names are not included in our data.

We can also see the [baby boomer period](#) quite clearly in the plot above.

Inferring Sex From Name

Let's use this dataset to estimate the number of females and males in our class. As with our class roster, we begin by lowercasing the names:

```

babynames['Name'] = babynames['Name'].str.lower()
babynames

```

1.4 What's in a Name?

	Name	Sex	Count	Year
0	mary	F	9217	1884
1	anna	F	3860	1884
2	emma	F	2587	1884
...
2081	verna	M	5	1883
2082	winnie	M	5	1883
2083	winthrop	M	5	1883

1891894 rows × 4 columns

Then, we count up how many male and female babies were born in total for each name:

```
sex_counts = pd.pivot_table(babynames, index='Name',
columns='Sex', values='Count',
aggfunc='sum', fill_value=0.,
margins=True)
sex_counts
```

Sex	F	M	All
Name			
aaban	0	96	96
aabha	35	0	35
aabid	0	10	10
...
zyyon	0	6	6
zzyzx	0	5	5
All	170639571	173894326	344533897

96175 rows × 3 columns

To decide whether a name is male or female, we can compute the proportion of times the name was given to a female baby.

```
prop_female = sex_counts['F'] / sex_counts['All']
sex_counts['prop_female'] = prop_female
sex_counts
```

Sex	F	M	All	prop_female
Name				
aaban	0	96	96	0.000000
aabha	35	0	35	1.000000
aabid	0	10	10	0.000000
...
zyyon	0	6	6	0.000000
zzyzx	0	5	5	0.000000
All	170639571	173894326	344533897	0.495277

96175 rows × 4 columns

We can then define a function that looks up the proportion of female names given a name.

```
def sex_from_name(name):
    if name in sex_counts.index:
        prop = sex_counts.loc[name, 'prop_female']
        return 'F' if prop > 0.5 else 'M'
    else:
        return None
sex_from_name('sam')
```

'M'

Try typing some names in this box to see whether the function outputs what you expect:

```
interact(sex_from_name, name='sam');
```

Show Widget

We mark each name in our class roster with its most likely sex.

```
students['sex'] = students['Name'].apply(sex_from_name)
students
```

	Name	Role	sex
0	keeley	Student	F
1	john	Student	M
2	bryan	Student	M
...
276	ernesto	Waitlist Student	M
277	athan	Waitlist Student	M
278	michael	Waitlist Student	M

279 rows × 3 columns

Now it is easy to estimate how many male and female students we have:

```
students['sex'].value_counts()
```

```
M    144
F    92
Name: sex, dtype: int64
```

Inferring Age From Name

We can proceed in a similar way to estimate the age distribution of the class, mapping each name to its average age in the dataset.

```

def avg_year(group):
    return np.average(group['Year'], weights=group['Count'])

avg_years = (
    babynames
    .groupby('Name')
    .apply(avg_year)
    .rename('avg_year')
    .to_frame()
)
avg_years

```

	avg_year
Name	
aaban	2012.572917
aabha	2013.714286
aabid	2009.500000
...	...
zyyanna	2010.000000
zyyon	2014.000000
zzyzx	2010.000000

96174 rows × 1 columns

```

def year_from_name(name):
    return (avg_years.loc[name, 'avg_year']
            if name in avg_years.index
            else None)

# Generate input box for you to try some names out:
interact(year_from_name, name='fernando');

```

Show Widget

```

students['year'] = students['Name'].apply(year_from_name)
students

```

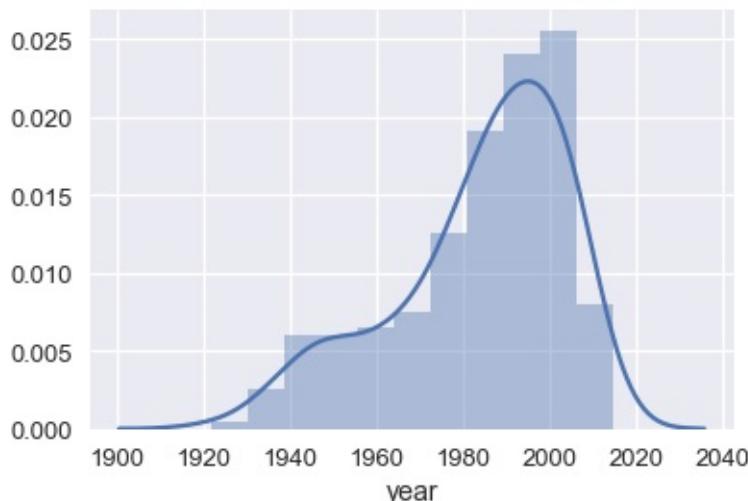
1.4 What's in a Name?

	Name	Role	sex	year
0	keeley	Student	F	1998.147952
1	john	Student	M	1951.084937
2	bryan	Student	M	1983.565113
...
276	ernesto	Waitlist Student	M	1981.439873
277	athan	Waitlist Student	M	2004.397863
278	michael	Waitlist Student	M	1971.179231

279 rows × 4 columns

Then, it is easy to plot the distribution of years:

```
sns.distplot(students['year'].dropna());
```



To compute the average year:

```
students['year'].mean()
```

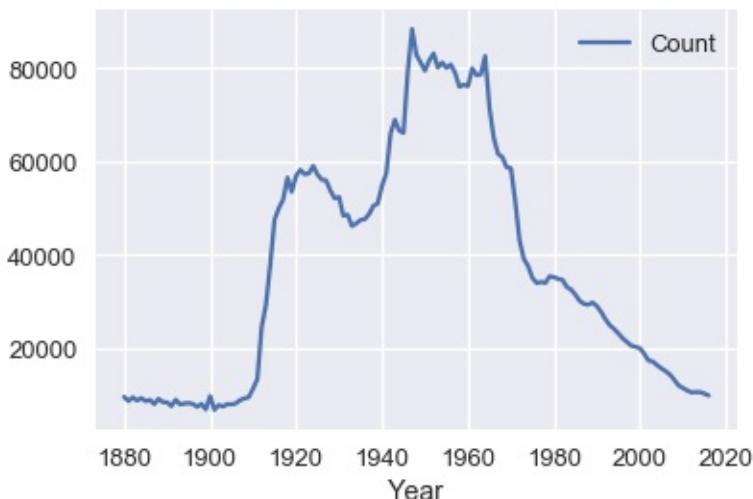
```
1983.846741800525
```

Which makes it appear like the average student is 35 years old. This is a course for college undergrads, so we were expecting an average age of around 20. Why might our estimate be so far off?

As data scientists, we often run into results that don't agree with our expectations and must make judgement calls about whether our results were caused by our data, our processes, or incorrect hypotheses. It is impossible to define rules that apply to every situation. Instead, we will give you the tools to re-examine each step of your data analysis and show you how to use them.

In this case, the most likely possibility for our unexpected result is that most names are old. For example, the name John was quite popular throughout the history recorded in our data, which means that we'd likely guess that a person named John is born around 1950. We can confirm this by looking at the data:

```
names = babynames.set_index('Name').sort_values('Year')
john = names.loc['john']
john[john['Sex'] == 'M'].plot('Year', 'Count');
```



If we happen to believe that no one in our class is over the age of 40 or under the age of 10 (we might find out by looking at our classroom during lecture) we can incorporate this into our analysis by only examining the data between 1978 and 2008. We will soon discuss data manipulation and you may revisit this analysis to find out if incorporating this prior belief gives a more sensible result.

Data Generation

Data science would hardly be a science without data. It is thus of utmost importance that we begin any data analysis by understanding how our data were generated.

In this chapter, we will discuss **data provenance**. Although the term data provenance typically refers to the entire history of data and where it has moved over time, we will use the term in this textbook to refer to the process that generated our data. Many have drawn premature conclusions because they were not careful enough in understanding their data; we will discuss an example of this to justify the importance of probability sampling in conducting data science.

Dewey Defeats Truman

In the 1948 US Presidential election, New York Governor Thomas Dewey ran against the incumbent Harry Truman. As usual, a number of polling agencies conducted polls of voters in order to predict which candidate was more likely to win the election.

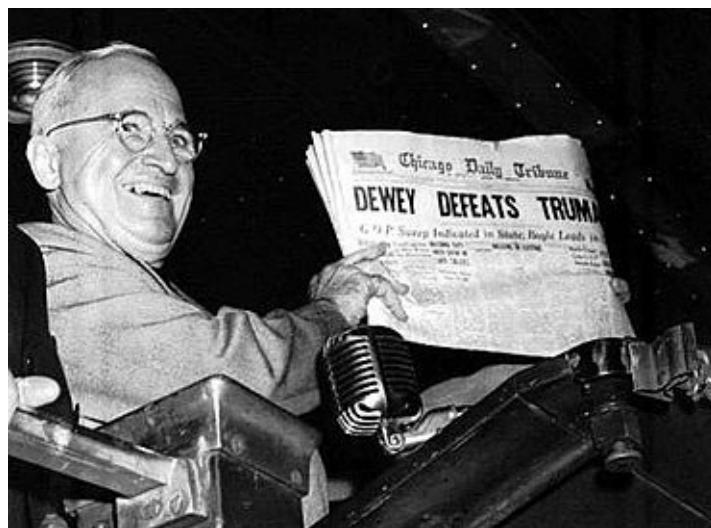
1936: A Previous Polling Catastrophe

In 1936, three elections prior to 1948, the *Literary Digest* infamously predicted a landslide defeat for Franklin Delano Roosevelt. To make this claim, the magazine polled a sample of over 2 million people based on telephone and car registrations. As you may know, this sampling scheme suffers from sampling bias: those with telephones and cars tend to be wealthier than those without. In this case, the sampling bias was so great that the *Literary Digest* thought Roosevelt would only receive 43% of the popular vote when he ended up with 61% of the popular vote, a difference of almost 20% and the largest error ever made by a major poll. The *Literary Digest* went out of business soon after¹.

1948: The Gallup Poll

Determined to learn from past mistakes, the Gallup Poll used a method called *quota sampling* to predict the results of the 1948 election. In their sampling scheme, each interviewer polled a set number of people from each demographic class. For example, the interviews were required to interview both males and females from different ages, ethnicities, and income levels to match the demographics in the US Census. This ensured that the poll would not leave out important subgroups of the voting population.

Using this method, the Gallup Poll predicted that Thomas Dewey would earn 5% more of the popular vote than Harry Truman would. This difference was significant enough that the *Chicago Tribune* famously printed the headline "Dewey Defeats Truman":



As we know now, Truman ended up winning the election. In fact, he won with 5% more of the popular vote than Dewey! What went wrong with the Gallup Poll?

The Problem With Quota Sampling

Although quota sampling did help pollsters reduce sampling bias, it introduced bias in another way. The Gallup Poll told its interviewers that as long as they fulfilled their quotas they could interview whomever they wished. Here's one possible explanation for why the interviewers ended up polling a disproportionate number of Republicans: at the time, Republicans were on average wealthier and more likely to live in nicer neighborhoods, making them easier to interview. This observation is supported by the fact that the Gallup Poll predicted 2-6% more Republican votes than the actual results for the 3 elections prior.

These examples highlight the importance of understanding sampling bias as much as possible during the data collection process. Both *Literary Digest* and Gallup Poll made the mistake of assuming their methods were unbiased when their sampling schemes were based on human judgement all along.

We now rely on **probability sampling**, a family of sampling methods that assigns precise probabilities to the appearance of each sample, to reduce bias as much as possible in our data collection process.

Big Data?

In the age of Big Data, we are tempted to deal with bias by collecting more data. After all, we know that a census will give us perfect estimates; shouldn't a very large sample give almost perfect estimates regardless of the sampling technique?

We will return to this question after discussing probability sampling methods to compare the two approaches.

1. <https://www.qualtrics.com/blog/the-1936-election-a-polling-catastrophe/> ↵

Probability Sampling

Unlike convenience sampling, probability sampling allows us to assign a precise probability to the event that we draw a particular sample. We will begin by reviewing simple random samples from Data 8, then introduce two alternative methods of probability sampling: cluster sampling and stratified sampling.

Suppose we have a population of 6 individuals. We've given each individual a different letter from A-F.

Simple Random Sample (SRS)

To take a simple random sample of size 2 from this population, we can write each letter A-F on a single index card, place all the cards into a hat, mix the cards well, and draw 2 cards without looking. That is, a SRS is sampling uniformly at random without replacement.

Here are all possible samples of size 2:

$$\begin{matrix} AB & BC & CD & DE & EF \\ AC&BD&CE&DF& \\ AD&BE&CF& \\ AE&BF& \\ AF& \end{matrix}$$

There are 15 possible samples of size 2 from our population of 6. Another way to count the number of possible samples is:

$$\binom{6}{2} = \frac{6!}{2!4!} = 15$$

Since in a SRS we sample uniformly at random, each of these 15 samples are equally likely to be chosen:

$$P(AB) = P(CD) = \dots = P(DF) = \frac{1}{15}$$

We can also use this chance mechanism to answer other questions about the composition of the sample. For example:

$$P(A \text{ in sample}) = \frac{5}{15} = \frac{1}{3}$$

Since 5 out of 15 of the possible samples listed above contain A.

By symmetry, we can say:

$$P(A \text{ in sample}) = P(F \text{ in sample}) = \frac{1}{3}$$

Another way of computing $P(A \text{ in sample})$ is to recognize that for A to be in the sample, we either need to draw it as the first marble or as the second marble.

$$\begin{aligned} P(A \text{ in sample}) &= P(\text{A is first or A is second}) \\ &= P(\text{A is first}) + P(\text{A is second}) \\ &= \frac{1}{6} \cdot \frac{5}{5} + \frac{5}{6} \cdot \frac{1}{5} \\ &= \frac{1}{3} \end{aligned}$$

Cluster Sampling

In cluster sampling, we divide the population into clusters. Then, we use SRS to select clusters at random instead of individuals.

As an example, suppose we take our population of 6 individuals and we pair each of them up: (A,B) (C,D) (E,F) to form 3 clusters of 2 individuals. Then, we use SRS to select one cluster to produce a sample of size 2.

As before, we can compute the probability that A is in our sample:

$$P(A \text{ in sample}) = P(AB \text{ drawn}) = \frac{1}{3}$$

Similarly the probability that any particular person appears in our sample is $\frac{1}{3}$. Note that this is the same as our SRS. However, we see differences when we look at the samples themselves. For example, in a SRS the chance of getting AB is the same as the chance of getting AC: $\frac{1}{15}$. However, with this cluster sampling scheme:

$$P(AB) = \frac{1}{3} \quad P(AC) = 0$$

Since A and C can never appear in the same sample if we only select one cluster.

Cluster sampling is still probability sampling since we can assign a probability to each potential sample. However, the resulting probabilities are different than using a SRS depending on how the population is clustered.

Why use cluster sampling? Cluster sampling is most useful because it makes sample collection easier. For example, it is much easier to poll towns of 100 people each than to poll thousands of people distributed across the entire US. This is the reason why many polling agencies today use forms of cluster sampling to conduct surveys.

The main downside of cluster sampling is that it tends to produce greater variation in estimation. This typically means that we take larger samples when using cluster sampling. Note that the reality is much more complicated than this, but we will leave the details to a future course on sampling techniques.

Stratified Sampling

In stratified sampling, we divide the population into strata, and then produce one simple random sample per strata. In both cluster sampling and stratified sampling we split the population into groups; in cluster sampling we use a single SRS to select groups whereas in

stratified sampling we use multiple SRS's, one for each group.

We can divide our population of 6 individuals into the following strata:

$$\begin{matrix} \text{Strata 1: } & \{A, B, C, D\} \\ \text{Strata 2: } & \{E, F\} \end{matrix}$$

We use an SRS to select one individual from each strata to produce a sample of size \$2\$.

This gives us the following possible samples:

$$\begin{aligned} (A, E) &\quad (A, F) & (B, E) &\quad (B, F) & (C, E) &\quad (C, F) \\ (D, E) &\quad (D, F) \end{aligned}$$

Again, we can compute the probability that \$A\$ is in our sample:

$$\begin{aligned} P(A \text{ in sample}) &= P(A \text{ selected from Strata 1}) \\ &= \frac{1}{4} \end{aligned}$$

However:

$$P(AB) = 0$$

since \$A\$ and \$B\$ cannot appear in the same sample.

Like cluster sampling, stratified sampling is also a probability sampling method that produces different probabilities depending on the stratification of the population. Note that like this example, the strata do not have to be the same size. For example, we can stratify the US by occupation, then take samples from each strata of size proportional to the distribution of occupations in the US — if only 0.01% of people in the US are statisticians, we can ensure that 0.01% of our sample will be composed of statisticians. A simple random sample might miss the poor statisticians altogether!

As you may have figured out, stratified sampling can perhaps be called the proper way to conduct quota sampling. It allows the researcher to ensure that subgroups of the population are well-represented in the sample without using human judgement to select the individuals in the sample. This can often result in less variation in estimation. However, stratified sampling is sometimes more difficult to accomplish because we sometimes don't know how large each strata is. In the previous example we have the advantage of the US census, but other times we are not so fortunate.

Why Probability Sampling?

As we have seen in Data 8, probability sampling enables us to quantify our uncertainty about an estimation or prediction. It is only through this precision that we can conduct inference and hypothesis testing. Be wary when anyone gives you p-values or confidence levels without a proper explanation of their sampling techniques.

Now that we understand probability sampling let us see how the humble SRS compares against "big data".

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [SRS vs. "Big Data"](#)
- [Takeaways](#)

SRS vs. "Big Data"

As we have previously mentioned, it is tempting to do away with our long-winded bias concerns by using huge amounts of data. It is true that collecting a census will by definition produce unbiased estimations. Perhaps we don't have to worry about bias if we just collect tons of data.

Suppose we are pollsters in 2012 trying to predict the popular vote of the US presidential election, where Barack Obama ran against Mitt Romney. Since we know the exact output of the popular vote, we can compare the predictions of a SRS to the predictions of a large non-random dataset, often called *administrative datasets* since they are often collected as part of some administrative work.

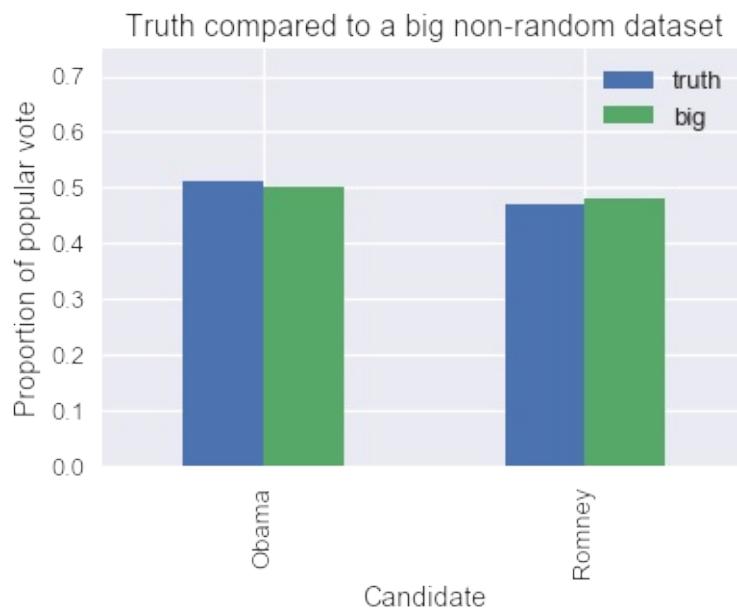
We will compare a SRS of size 400 to a non-random sample of size 60,000,000. Our non-random sample is nearly 150,000 times larger than our SRS! Since there were about 120,000,000 voters in 2012, we can think of our non-random sample as a survey where half of all voters in the US responded (no actual poll has ever surveyed more than 10,000,000 voters).

Here's a plot comparing the proportions of the non-random sample to the true proportions. The bars labeled `truth` show the true proportions of votes that each candidate received. The bars labeled `big` show the proportions from our dataset of 60,000,000 voters.

```

pd.DataFrame({
    'truth': [obama_true, romney_true],
    'big': [obama_big, romney_big],
}, index=['Obama', 'Romney'], columns=['truth',
'big']).plot.bar()
plt.title('Truth compared to a big non-random dataset')
plt.xlabel('Candidate')
plt.ylabel('Proportion of popular vote')
plt.ylim(0, 0.75)
None

```



We can see that our large dataset is just a bit biased towards the Republican candidate Romney just as the Gallup Poll was in 1948. Still, this dataset could give us accurate predictions. To check, we can simulate taking simple random samples of size 400 from the population and large non-random samples of size 60,000,000. We will compute the proportion of votes for Obama in each sample and plot the distribution of proportions.

```
srs_size = 400
big_size = 60000000
replications = 10000

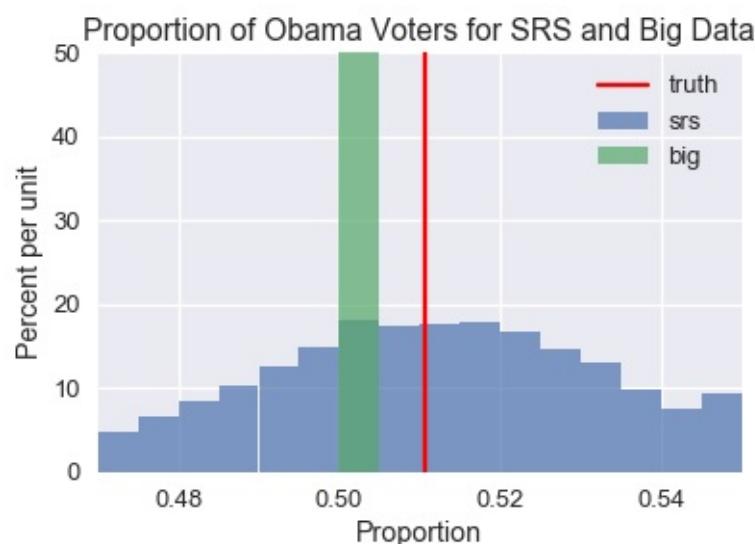
def resample(size, prop, replications):
    return np.random.binomial(n=size, p=prop, size=replications)
/ size

srs_simulations = resample(srs_size, obama_true, replications)
big_simulations = resample(big_size, obama_big, replications)
```

Now, we will plot the simulation results and overlay a red line indicating the true proportion of voters that voted for Obama.

```
bins = bins=np.arange(0.47, 0.55, 0.005)
plt.hist(srs_simulations, bins=bins, alpha=0.7, normed=True,
label='srs')
plt.hist(big_simulations, bins=bins, alpha=0.7, normed=True,
label='big')

plt.title('Proportion of Obama Voters for SRS and Big Data')
plt.xlabel('Proportion')
plt.ylabel('Percent per unit')
plt.xlim(0.47, 0.55)
plt.ylim(0, 50)
plt.axvline(x=obama_true, color='r', label='truth')
plt.legend()
None
```



As you can see, the SRS distribution is spread out but centered around the true population proportion of Obama voters. The distribution created by the large non-random sample, on the other hand, is very narrow but not a single simulated sample produces the true population proportion. If we attempt to create confidence intervals using the non-random sample, none of them will contain the true population proportion. To make matters worse, the confidence interval will be extremely narrow because the sample is so large. We will be very sure of an ultimately incorrect estimation.

In fact, when our sampling method is biased our estimations will become **worse** as we collect more data since we will be more certain about an incorrect result, only becoming more accurate when our dataset is almost a census of our population. **The quality of the data matters much more than its size.**

Takeaways

Before accepting the results of a data analysis, it pays to carefully inspect the quality of the data. In particular, we must ask the following questions:

1. Is the data a census (does it include the entire population of interest)? If so, we can just compute properties of the population directly without having to use inference.
2. If the data is a sample, how was the sample collected? To properly conduct inference, the sample should have been collected according to a probability sampling method.
3. What changes were made to the data before producing results? Do any of these changes affect the quality of the data?

For more details on the comparison between random and large non-random samples, we suggest watching [this lecture by the statistician Xiao-Li Meng](#).

Working with Tabular Data

Tabular data, like the datasets we have worked with in Data 8, are one of the most common and useful forms of data for analysis. We introduce tabular data manipulation using `pandas`, the standard Python library for working with tabular data. Although `pandas`'s syntax is more challenging to use than the `datascience` package used in Data 8, `pandas` provides significant performance improvements and is the current tool of choice in both industry and academia for working with tabular data.

It is more important that you understand the types of useful operations on data than the exact details of `pandas` syntax. For example, knowing when to use a group or a join is more useful than knowing how to call the `pandas` function to group data. It is relatively easy to look up the function you need once you know the right operation to use. All of the table manipulations in this chapter will also appear again in a new syntax when we cover SQL, so it will help you to understand them now.

Because we will cover only the most important `pandas` functions in this textbook, you should bookmark the [pandas documentation](#) for reference when you conduct your own data analyses.

Show Widgets [Open on DataHub](#)

Table of Contents

- [Getting Started](#)
- [Indexes, Slicing, and Sorting](#)
 - [Breaking the Problem Down](#)
 - [Slicing using `.loc`](#)
 - [Slicing rows using a predicate](#)
 - [Sorting Rows](#)
- [In Conclusion](#)

Getting Started

In each section of this chapter we will work with the Baby Names dataset from Chapter 1. We will pose a question, break the question down into high-level steps, then translate each step into Python code using `pandas` `DataFrames`. We begin by importing `pandas`:

```
# pd is a common shorthand for pandas
import pandas as pd
```

Now we can read in the data using `pd.read_csv` ([docs](#)).

```
baby = pd.read_csv('babynames.csv')
baby
```

	Name	Sex	Count	Year
0	Mary	F	9217	1884
1	Anna	F	3860	1884
2	Emma	F	2587	1884
...
1891891	Verna	M	5	1883
1891892	Winnie	M	5	1883
1891893	Winthrop	M	5	1883

1891894 rows × 4 columns

Note that for the code above to work, the `babynames.csv` file must be located in the same directory as this notebook. We can check what files are in the current folder by running `ls` in a notebook cell:

```
ls
```

```
babynames.csv
```

```
indexes_slicing_sorting.ipynb
```

When we use `pandas` to read in data, we get a DataFrame. A DataFrame is a tabular data structure where each column is labeled (in this case 'Name', 'Sex', 'Count', 'Year') and each row is labeled (in this case 0, 1, 2, ..., 1891893). The Table introduced in Data 8, however, only has labeled columns.

The labels of a DataFrame are called the *indexes* of the DataFrame and make many data manipulations easier.

Indexes, Slicing, and Sorting

Let's use `pandas` to answer the following question:

What were the five most popular baby names in 2016?

Breaking the Problem Down

We can decompose this question into the following simpler table manipulations:

1. Slice out the rows for the year 2016.
2. Sort the rows in descending order by Count.

Now, we can express these steps in `pandas`.

Slicing using `.loc`

To select subsets of a DataFrame, we use the `.loc` slicing syntax. The first argument is the label of the row and the second is the label of the column:

```
baby
```

	Name	Sex	Count	Year
0	Mary	F	9217	1884
1	Anna	F	3860	1884
2	Emma	F	2587	1884
...
1891891	Verna	M	5	1883
1891892	Winnie	M	5	1883
1891893	Winthrop	M	5	1883

1891894 rows × 4 columns

```
baby.loc[1, 'Name'] # Row labeled 1, Column labeled 'Name'
```

```
'Anna'
```

To slice out multiple rows or columns, we can use `:`. Note that `.loc` slicing is inclusive, unlike Python's slicing.

```
# Get rows 1 through 5, columns Name through Count inclusive
baby.loc[1:5, 'Name':'Count']
```

	Name	Sex	Count
1	Anna	F	3860
2	Emma	F	2587
3	Elizabeth	F	2549
4	Minnie	F	2243
5	Margaret	F	2142

We will often want a single column from a DataFrame:

```
baby.loc[:, 'Year']
```

```
0          1884
1          1884
2          1884
...
1891891    1883
1891892    1883
1891893    1883
Name: Year, Length: 1891894, dtype: int64
```

Note that when we select a single column, we get a `pandas Series`. A Series is like a one-dimensional NumPy array since we can perform arithmetic on all the elements at once.

```
baby.loc[:, 'Year'] * 2
```

```
0          3768
1          3768
2          3768
...
1891891    3766
1891892    3766
1891893    3766
Name: Year, Length: 1891894, dtype: int64
```

To select out specific columns, we can pass a list into the `.loc` slice:

```
# This is a DataFrame again
baby.loc[:, ['Name', 'Year']]
```

	Name	Year
0	Mary	1884
1	Anna	1884
2	Emma	1884
...
1891891	Verna	1883
1891892	Winnie	1883
1891893	Winthrop	1883

1891894 rows × 2 columns

Selecting columns is common, so there's a shorthand.

```
# Shorthand for baby.loc[:, 'Name']
baby['Name']
```

```
0      Mary
1      Anna
2      Emma
...
1891891    Verna
1891892    Winnie
1891893    Winthrop
Name: Name, Length: 1891894, dtype: object
```

```
# Shorthand for baby.loc[:, ['Name', 'Count']]
baby[['Name', 'Count']]
```

	Name	Count
0	Mary	9217
1	Anna	3860
2	Emma	2587
...
1891891	Verna	5
1891892	Winnie	5
1891893	Winthrop	5

1891894 rows × 2 columns

Slicing rows using a predicate

To slice out the rows with year 2016, we will first create a Series containing `True` for each row we want to keep and `False` for each row we want to drop. This is simple because math and boolean operators on Series are applied to each element in the Series.

```
# Series of years
baby['Year']
```

```
0      1884
1      1884
2      1884
...
1891891    1883
1891892    1883
1891893    1883
Name: Year, Length: 1891894, dtype: int64
```

```
# Compare each year with 2016
baby['Year'] == 2016
```

```

0      False
1      False
2      False
...
1891891  False
1891892  False
1891893  False
Name: Year, Length: 1891894, dtype: bool

```

Once we have this Series of `True` and `False`, we can pass it into `.loc`.

```

# We are slicing rows, so the boolean Series goes in the first
# argument to .loc
baby_2016 = baby.loc[baby['Year'] == 2016, :]
baby_2016

```

	Name	Sex	Count	Year
1850880	Emma	F	19414	2016
1850881	Olivia	F	19246	2016
1850882	Ava	F	16237	2016
...
1883745	Zyahir	M	5	2016
1883746	Zyel	M	5	2016
1883747	Zylyn	M	5	2016

32868 rows × 4 columns

Sorting Rows

The next step is to sort the rows in descending order by 'Count'. We can use the `sort_values()` function.

```

sorted_2016 = baby_2016.sort_values('Count', ascending=False)
sorted_2016

```

	Name	Sex	Count	Year
1850880	Emma	F	19414	2016
1850881	Olivia	F	19246	2016
1869637	Noah	M	19015	2016
...
1868752	Mikaelyn	F	5	2016
1868751	Miette	F	5	2016
1883747	Zylyn	M	5	2016

32868 rows × 4 columns

Finally, we will use `.iloc` to slice out the first five rows of the DataFrame. `.iloc` works like `.loc` but takes in numerical indices instead of labels. It does not include the right endpoint in its slices, like Python's list slicing.

```
# Get the value in the zeroth row, zeroth column
sorted_2016.iloc[0, 0]
```

```
'Emma'
```

```
# Get the first five rows
sorted_2016.iloc[0:5]
```

	Name	Sex	Count	Year
1850880	Emma	F	19414	2016
1850881	Olivia	F	19246	2016
1869637	Noah	M	19015	2016
1869638	Liam	M	18138	2016
1850882	Ava	F	16237	2016

In Conclusion

We now have the five most popular baby names in 2016 and learned to express the following operations in `pandas`:

Operation	pandas
Read a CSV file	<code>pd.read_csv()</code>
Slicing using labels or indices	<code>.loc</code> and <code>.iloc</code>
Slicing rows using a predicate	Use a boolean-valued Series in <code>.loc</code>
Sorting rows	<code>.sort_values()</code>

Show Widgets [Open on DataHub](#)

Table of Contents

- [Grouping and Pivoting](#)
 - [Breaking the Problem Down](#)
 - [Grouping](#)
 - [Grouping on Multiple Columns](#)
 - [Pivoting](#)
- [In Conclusion](#)

Grouping and Pivoting

In this section, we will answer the question:

What were the most popular male and female names in each year?

Here's the Baby Names dataset once again:

```
baby = pd.read_csv('babynames.csv')
baby.head()
# the .head() method outputs the first five rows of the
DataFrame
```

	Name	Sex	Count	Year
0	Mary	F	9217	1884
1	Anna	F	3860	1884
2	Emma	F	2587	1884
3	Elizabeth	F	2549	1884
4	Minnie	F	2243	1884

Breaking the Problem Down

We should first notice that the question in the previous section has similarities to this one; the question in the previous section restricts names to babies born in 2016 whereas this question asks for names in all years.

We once again decompose this problem into simpler table manipulations.

1. Group the `baby` DataFrame by 'Year' and 'Sex'.
2. For each group, compute the most popular name.

Recognizing which operation is needed for each problem is sometimes tricky. Usually, a convoluted series of steps will signal to you that there might be a simpler way to express what you want. If we didn't immediately recognize that we needed to group, for example, we might write steps like the following:

1. Loop through each unique year.
2. For each year, loop through each unique sex.
3. For each unique year and sex, find the most common name.

There is almost always a better alternative to looping over a `pandas` DataFrame. **In particular, looping over unique values of a DataFrame should usually be replaced with a group.**

Grouping

To group in `pandas`, we use the `.groupby()` method.

```
baby.groupby('Year')
```

```
<pandas.core.groupby.DataFrameGroupBy object at 0x1a14e21f60>
```

`.groupby()` returns a strange-looking `DataFrameGroupBy` object. We can call `.agg()` on this object with an aggregation function in order to get a familiar output:

```
# The aggregation function takes in a series of values for each
# group
# and outputs a single value
def length(series):
    return len(series)

# Count up number of values for each year. This is equivalent to
# counting the number of rows where each year appears.
baby.groupby('Year').agg(length)
```

	Name	Sex	Count
Year			
1880	2000	2000	2000
1881	1935	1935	1935
1882	2127	2127	2127
...
2014	33206	33206	33206
2015	33063	33063	33063
2016	32868	32868	32868

137 rows × 3 columns

You might notice that the `length` function simply calls the `len` function, so we can simplify the code above.

```
baby.groupby('Year').agg(len)
```

	Name	Sex	Count
Year			
1880	2000	2000	2000
1881	1935	1935	1935
1882	2127	2127	2127
...
2014	33206	33206	33206
2015	33063	33063	33063
2016	32868	32868	32868

137 rows × 3 columns

The aggregation is applied to each column of the DataFrame, producing redundant information. We can restrict the output columns by slicing before grouping.

```

year_rows = baby[['Year', 'Count']].groupby('Year').agg(len)
year_rows

# A further shorthand to accomplish the same result:
#
# year_counts = baby[['Year', 'Count']].groupby('Year').count()
#
# pandas has shorthands for common aggregation functions,
# including
# count, sum, and mean.

```

	Count
Year	
1880	2000
1881	1935
1882	2127
...	...
2014	33206
2015	33063
2016	32868

137 rows × 1 columns

Note that the index of the resulting DataFrame now contains the unique years, so we can slice subsets of years using `.loc` as before:

```

# Every twentieth year starting at 1880
year_rows.loc[1880:2016:20, :]

```

	Count
Year	
1880	2000
1900	3730
1920	10755
1940	8961
1960	11924
1980	19440
2000	29764

Grouping on Multiple Columns

As we've seen in Data 8, we can group on multiple columns to get groups based on unique pairs of values. To do this, pass in a list of column labels into `.groupby()`.

```
grouped_counts = baby.groupby(['Year', 'Sex']).sum()
grouped_counts
```

		Count
Year	Sex	
1880	F	90992
	M	110491
1881	F	91953
...
2015	M	1907211
2016	F	1756647
	M	1880674

274 rows × 1 columns

The code above computes the total number of babies born for each year and sex. Let's now use grouping by multiple columns to compute the most popular names for each year and sex. Since the data are already sorted in descending order of Count for each year and sex, we can define an aggregation function that returns the first value in each series. (If the data weren't sorted, we can call `sort_values() first.`)

```
# The most popular name is simply the first one that appears in
# the series
def most_popular(series):
    return series.iloc[0]

baby_pop = baby.groupby(['Year', 'Sex']).agg(most_popular)
baby_pop
```

		Name	Count
Year	Sex		
1880	F	Mary	7065
	M	John	9655
1881	F	Mary	6919
...
2015	M	Noah	19594
2016	F	Emma	19414
	M	Noah	19015

274 rows × 2 columns

Notice that grouping by multiple columns results in multiple labels for each row. This is called a "multilevel index" and is tricky to work with. The important thing to know is that `.loc` takes in a tuple for the row index instead of a single value:

```
baby_pop.loc[(2000, 'F'), 'Name']
```

```
'Emily'
```

But `.iloc` behaves the same as usual since it uses indices instead of labels:

```
baby_pop.iloc[10:15, :]
```

		Name	Count
Year	Sex		
1885	F	Mary	9128
	M	John	8756
1886	F	Mary	9889
	M	John	9026
1887	F	Mary	9888

Pivoting

If you group by two columns, you can often use pivot to present your data in a more convenient format. Using a pivot lets you use one set of grouped labels as the columns of the resulting table.

To pivot, use the `pd.pivot_table()` function.

```
pd.pivot_table(baby,
               index='Year',           # Index for rows
               columns='Sex',          # Columns
               values='Name',          # Values in table
               aggfunc=most_popular) # Aggregation function
```

Sex	F	M
Year		
1880	Mary	John
1881	Mary	John
1882	Mary	John
...
2014	Emma	Noah
2015	Emma	Noah
2016	Emma	Noah

137 rows × 2 columns

Compare this result to the `baby_pop` table that we computed using `.groupby()`. We can see that the `Sex` index in `baby_pop` became the columns of the pivot table.

```
baby_pop
```

		Name	Count
Year	Sex		
1880	F	Mary	7065
	M	John	9655
1881	F	Mary	6919
...
2015	M	Noah	19594
2016	F	Emma	19414
	M	Noah	19015

274 rows × 2 columns

In Conclusion

We now have the most popular baby names for each sex and year in our dataset and learned to express the following operations in pandas :

Operation	pandas
Group	<code>df.groupby(label)</code>
Group by multiple columns	<code>df.groupby([label1, label2])</code>
Group and aggregate	<code>df.groupby(label).agg(func)</code>
Pivot	<code>pd.pivot_table()</code>

Show Widgets [Open on DataHub](#)

Table of Contents

- [Apply, Strings, and Plotting](#)
 - [Breaking the Problem Down](#)
 - [Apply](#)
 - [String Manipulation](#)
 - [Grouping](#)
 - [Plotting](#)
- [In Conclusion](#)

Apply, Strings, and Plotting

In this section, we will answer the question:

Can we use the last letter of a name to predict the sex of the baby?

Here's the Baby Names dataset once again:

```
baby = pd.read_csv('babynames.csv')
baby.head()
# the .head() method outputs the first five rows of the
DataFrame
```

	Name	Sex	Count	Year
0	Mary	F	9217	1884
1	Anna	F	3860	1884
2	Emma	F	2587	1884
3	Elizabeth	F	2549	1884
4	Minnie	F	2243	1884

Breaking the Problem Down

Although there are many ways to see whether prediction is possible, we will use plotting in this section. We can decompose this question into two steps:

1. Compute the last letter of each name.

2. Group by the last letter and sex, aggregating on Count.
3. Plot the counts for each sex and letter.

Apply

pandas Series contain an `.apply()` method that takes in a function and applies it to each value in the Series.

```
names = baby['Name']
names.apply(len)
```

```
0          4
1          4
2          4
..
1891891    5
1891892    6
1891893    8
Name: Name, Length: 1891894, dtype: int64
```

To extract the last letter of each name, we can define our own function to pass into `.apply()`:

```
def last_letter(string):
    return string[-1]

names.apply(last_letter)
```

```
0          y
1          a
2          a
..
1891891    a
1891892    e
1891893    p
Name: Name, Length: 1891894, dtype: object
```

String Manipulation

Although `.apply()` is flexible, it is often faster to use the built-in string manipulation functions in `pandas` when dealing with text data.

`pandas` provides access to string manipulation functions using the `.str` attribute of Series.

```
names = baby['Name']
names.str.len()
```

```
0          4
1          4
2          4
...
1891891    5
1891892    6
1891893    8
Name: Name, Length: 1891894, dtype: int64
```

We can directly slice out the last letter of each name in a similar way.

```
names.str[-1]
```

```
0          y
1          a
2          a
...
1891891    a
1891892    e
1891893    p
Name: Name, Length: 1891894, dtype: object
```

We suggest looking at the docs for the full list of string methods ([link](#)).

We can now add this column of last letters to our `baby` DataFrame.

```
baby['Last'] = names.str[-1]
baby
```

	Name	Sex	Count	Year	Last
0	Mary	F	9217	1884	y
1	Anna	F	3860	1884	a
2	Emma	F	2587	1884	a
...
1891891	Verna	M	5	1883	a
1891892	Winnie	M	5	1883	e
1891893	Winthrop	M	5	1883	p

1891894 rows × 5 columns

Grouping

To compute the sex distribution for each last letter, we need to group by both Last and Sex.

```
# Shorthand for baby.groupby(['Last', 'Sex']).agg(np.sum)
baby.groupby(['Last', 'Sex']).sum()
```

		Count	Year
Last	Sex		
a	F	58079486	915565667
	M	1931630	53566324
b	F	17376	1092953
...
y	M	18569388	114394474
z	F	142023	4268028
	M	120123	9649274

52 rows × 2 columns

Notice that `Year` is also summed up since each non-grouped column is passed into the aggregation function. To avoid this, we can select out the desired columns before calling `.groupby()`.

```
# When lines get long, you can wrap the entire expression in
# parentheses
# and insert newlines before each method call
letter_dist = (
    baby[['Last', 'Sex', 'Count']]
    .groupby(['Last', 'Sex'])
    .sum()
)
letter_dist
```

		Count
Last	Sex	
a	F	58079486
	M	1931630
b	F	17376
...
y	M	18569388
z	F	142023
	M	120123

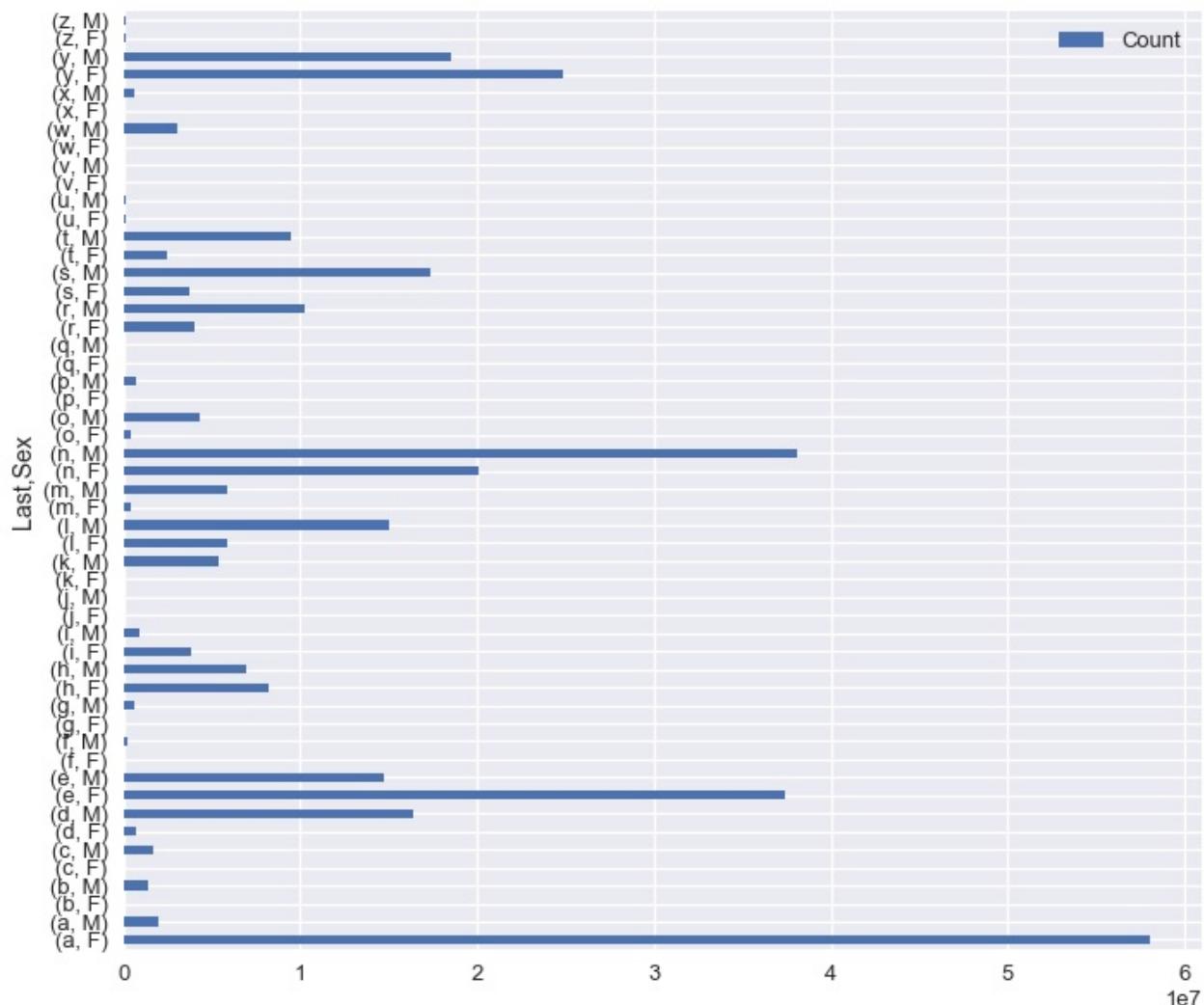
52 rows × 1 columns

Plotting

`pandas` provides built-in plotting functionality for most basic plots, including bar charts, histograms, line charts, and scatterplots. To make a plot from a DataFrame, use the `.plot` attribute:

```
# We use the figsize option to make the plot larger
letter_dist.plot.barch(figsize=(10, 10))
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a17af4780>
```



Although this plot shows the distribution of letters and sexes, the male and female bars are difficult to tell apart. By looking at the `pandas` docs on plotting ([link](#)) we learn that `pandas` plots one group of bars for row column in the DataFrame, showing one differently colored bar for each column. This means that a pivoted version of the `letter_dist` table will have the right format.

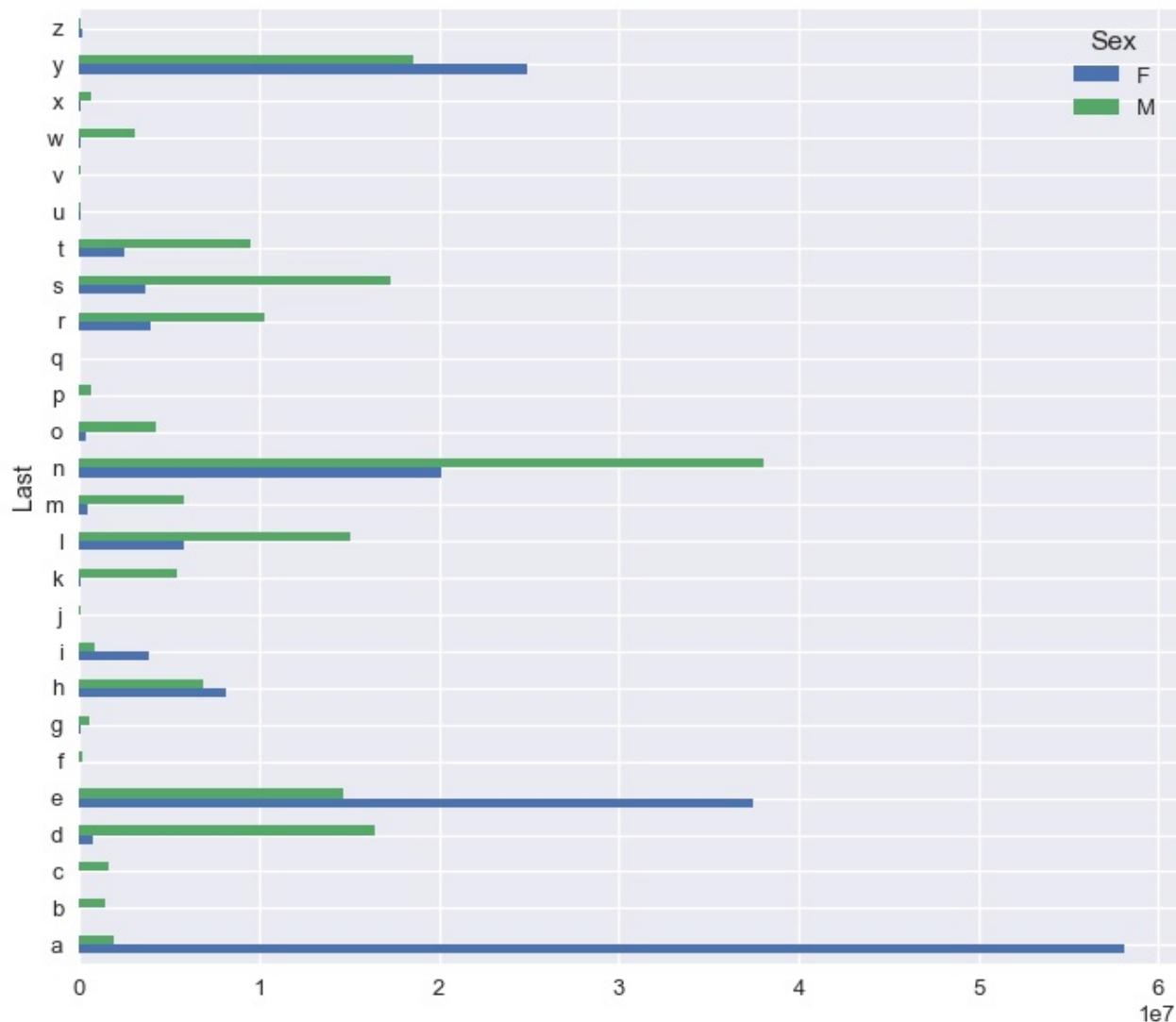
```
letter_pivot = pd.pivot_table(
    baby, index='Last', columns='Sex', values='Count',
    aggfunc='sum'
)
letter_pivot
```

Sex	F	M
Last		
a	58079486	1931630
b	17376	1435939
c	30262	1672407
...
x	37381	644092
y	24877638	18569388
z	142023	120123

26 rows × 2 columns

```
letter_pivot.plot.barh(figsize=(10, 10))
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a17c36978>
```



Notice that `pandas` conveniently generates a legend for us as well. However, this is still difficult to interpret. We plot the counts for each letter and sex which causes some bars to appear very long and others to be almost invisible. We should instead plot the proportion of male and female babies within each last letter.

```
total_for_each_letter = letter_pivot['F'] + letter_pivot['M']

letter_pivot['F prop'] = letter_pivot['F'] /
total_for_each_letter
letter_pivot['M prop'] = letter_pivot['M'] /
total_for_each_letter
letter_pivot
```

Sex	F	M	F prop	M prop
Last				
a	58079486	1931630	0.967812	0.032188
b	17376	1435939	0.011956	0.988044
c	30262	1672407	0.017773	0.982227
...
x	37381	644092	0.054853	0.945147
y	24877638	18569388	0.572597	0.427403
z	142023	120123	0.541771	0.458229

26 rows × 4 columns

```
(letter_pivot[['F prop', 'M prop']]  
 .sort_values('M prop') # Sorting orders the plotted bars  
 .plot.barh(figsize=(10, 10))  
)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a18194b70>
```



In Conclusion

We can see that almost all first names that end in 'p' are male and names that end in 'a' are female! In general, the difference between bar lengths for many letters implies that we can often make a good guess to a person's sex if we just know the last letter of their first name.

We've learned to express the following operations in `pandas` :

Operation	pandas
Applying a function elementwise	<code>series.apply(func)</code>
String manipulation	<code>series.str.func()</code>
Plotting	<code>df.plot.func()</code>

Data Cleaning

Data come in many formats and vary greatly in usefulness for analysis. Although we would prefer all our data to come in a tabular format with each value recorded consistently and accurately, in reality we must carefully check our data for potential issues that can eventually result in incorrect conclusions.

The term "data cleaning" refers to the process of combing through the data and deciding how to resolve inconsistencies and missing values. We will discuss common problems found in datasets and approaches to address them.

Data cleaning has inherent limitations. For example, no amount of data cleaning will fix a biased sampling process. Before embarking on the sometimes lengthy process of data cleaning, we must be confident that our data are collected accurately and with as little bias as possible. Only then can we investigate the data itself and use data cleaning to resolve issues in the data format or entry process.

We will introduce data cleaning techniques by working with City of Berkeley Police Department datasets.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Investigating Berkeley Police Data](#)
 - [Understanding the Data Generation](#)
- [Cleaning The Calls Dataset](#)
 - [Are there missing values?](#)
 - [Are there any missing values that were filled in?](#)
 - [Which parts of the data were entered by a human?](#)
 - [Final Touchups](#)

Investigating Berkeley Police Data

We will use the Berkeley Police Department's publicly available datasets to demonstrate data cleaning techniques. We have downloaded the [Calls for Service dataset](#) and [Stops dataset](#).

We can use the `ls` shell command with the `-lh` flags to see more details about the files:

```
!ls -lh data/
```

```
total 13936
-rw-r--r--@ 1 sam  staff  979K Aug 29 14:41 Berkeley_PD_-_Calls_for_Service.csv
-rw-r--r--@ 1 sam  staff   81B Aug 29 14:28 cvdow.csv
-rw-r--r--@ 1 sam  staff  5.8M Aug 29 14:41 stops.json
```

The command above shows the data files and their file sizes. This is especially useful because we now know the files are small enough to load into memory. As a rule of thumb, it is usually safe to load a file into memory that is around one fourth of the total memory capacity of the computer. For example, if a computer has 4GB of RAM we should be able to load a 1GB CSV file in `pandas`. To handle larger datasets we will need additional computational tools that we will cover later in this book.

Notice the use of the exclamation point before `ls`. This tells Jupyter that the next line of code is a shell command, not a Python expression. We can run any available shell command in Jupyter using `!`:

```
# The `wc` shell command shows us how many lines each file has.  
# We can see that the `stops.json` file has the most lines  
(29852).  
!wc -l data/*
```

```
16497 data/Berkeley_PD_-_Calls_for_Service.csv  
     8 data/cvdow.csv  
29852 data/stops.json  
46357 total
```

Understanding the Data Generation

We will state important questions you should ask of all datasets before data cleaning or processing. These questions are related to how the data were generated, so data cleaning will usually **not** be able to resolve issues that arise here.

What do the data contain? The website for the Calls for Service data states that the dataset describes "crime incidents (not criminal reports) within the last 180 days". Further reading reveals that "not all calls for police service are included (e.g. Animal Bite)".

The website for the Stops data states that the dataset contains data on all "vehicle detentions (including bicycles) and pedestrian detentions (up to five persons)" since January 26, 2015.

Are the data a census? This depends on our population of interest. For example, if we are interested in calls for service within the last 180 days for crime incidents then the Calls dataset is a census. However, if we are interested in calls for service within the last 10 years the dataset is clearly not a census. We can make similar statements about the Stops dataset since the data collection started on January 26, 2015.

If the data form a sample, is it a probability sample? If we are investigating a period of time that the data do not have entries for, the data do not form a probability sample since there is no randomness involved in the data collection process — we have all data for certain time periods but no data for others.

What limitations will this data have on our conclusions? Although we will ask this question at each step of our data processing, we can already see that our data impose important limitations. The most important limitation is that we cannot make unbiased estimations for time periods not covered by our datasets.

Cleaning The Calls Dataset

Let's now clean the Calls dataset. The `head` shell command prints the first five lines of the file.

```
!head data/Berkeley_PD_-_Calls_for_Service.csv
```

```
CASENO,OFFENSE,EVENTDT,EVENTTM,CVLEGEND,CVDOW,InDbDate,Block_Location,BLKADDR,City,State  
17091420,BURGLARY AUTO,07/23/2017 12:00:00 AM,06:00,BURGLARY - VEHICLE,0,08/29/2017 08:28:05 AM,"2500 LE CONTE AVE Berkeley, CA (37.876965, -122.260544)",2500 LE CONTE AVE,Berkeley,CA  
17020462,THEFT FROM PERSON,04/13/2017 12:00:00 AM,08:45,LARCENY,4,08/29/2017 08:28:00 AM,"2200 SHATTUCK AVE Berkeley, CA (37.869363, -122.268028)",2200 SHATTUCK AVE,Berkeley,CA  
17050275,BURGLARY AUTO,08/24/2017 12:00:00 AM,18:30,BURGLARY - VEHICLE,4,08/29/2017 08:28:06 AM,"200 UNIVERSITY AVE Berkeley, CA (37.865491, -122.310065)",200 UNIVERSITY AVE,Berkeley,CA
```

It appears to be a comma-separated values (CSV) file, though it's hard to tell whether the entire file is formatted properly. We can use `pd.read_csv` to read in the file as a DataFrame. If `pd.read_csv` errors, we will have to dig deeper and manually resolve formatting issues. Fortunately, `pd.read_csv` successfully returns a DataFrame:

```
calls = pd.read_csv('data/Berkeley_PD_-_Calls_for_Service.csv')  
calls
```

4.1 Cleaning the Calls Dataset

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	Block_Loc
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	2500 LE CC AVE\nBerke CA\n(37.876 -...
1	17020462	THEFT FROM PERSON	04/13/2017 12:00:00 AM	08:45	...	2200 SHATTUCK AVE\nBerke CA\n(37.869 -...
2	17050275	BURGLARY AUTO	08/24/2017 12:00:00 AM	18:30	...	200 UNIVERSIT AVE\nBerke CA\n(37.865 ...
...
5505	17018126	DISTURBANCE	04/01/2017 12:00:00 AM	12:22	...	1600 FAIRV ST\nBerkele CA\n(37.850 -1...
5506	17090665	THEFT MISD. (UNDER \$950)	04/01/2017 12:00:00 AM	12:00	...	2000 DELAWARE ST\nBerkele CA\n(37.874 -1...
5507	17049700	SEXUAL ASSAULT MISD.	08/22/2017 12:00:00 AM	20:02	...	2400 TELEGRAP AVE\nBerke CA\n(37.866 ...

5508 rows × 11 columns

We can define a function to show different slices of the data and then interact with it:

```

def df_interact(df):
    """
    Outputs sliders that show rows and columns of df
    """

    def peek(row=0, col=0):
        return df.iloc[row:row + 5, col:col + 6]
    interact(peek, row=(0, len(df), 5), col=(0, len(df.columns) - 6))
    print('({} rows, {} columns) total'.format(df.shape[0], df.shape[1]))

df_interact(calls)

```

Show Widget

(5508 rows, 11 columns) total

Based on the output above, the resulting DataFrame looks reasonably well-formed since the columns are properly named and the data in each column seems to be entered consistently. What data does each column contain? We can look at the dataset website:

Column	Description	Type
CASENO	Case Number	Number
OFFENSE	Offense Type	Plain Text
EVENTDT	Date Event Occurred	Date & Time
EVENTTM	Time Event Occurred	Plain Text
CVLEGEND	Description of Event	Plain Text
CVDOW	Day of Week Event Occurred	Number
InDbDate	Date dataset was updated in the portal	Date & Time
Block_Location	Block level address of event	Location
BLKADDR		Plain Text
City		Plain Text
State		Plain Text

On the surface the data looks easy to work with. However, before starting data analysis we must answer the following questions:

1. **Are there missing values in the dataset?** This question is important because missing

values can represent many different things. For example, missing addresses could mean that locations were removed to protect anonymity, or that some respondents chose not to answer a survey question, or that a recording device broke.

2. **Are there any missing values that were filled in (e.g. a 999 for unknown age or 12:00am for unknown date)?** These will clearly impact analysis if we ignore them.
3. **Which parts of the data were entered by a human?** As we will soon see, human-entered data is filled with inconsistencies and misspellings.

Although there are plenty more checks to go through, these three will suffice for many cases. See the [Quartz bad data guide](#) for a more complete list of checks.

Are there missing values?

This is a simple check in `pandas`:

```
# True if row contains at least one null value
null_rows = calls.isnull().any(axis=1)
calls=null_rows]
```

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	Block_Locat
116	17014831	BURGLARY AUTO	03/16/2017 12:00:00 AM	22:00	...	Berkeley, CA\n(37.8690 -122.270455)
478	17042511	BURGLARY AUTO	07/20/2017 12:00:00 AM	16:00	...	Berkeley, CA\n(37.8690 -122.270455)
486	17022572	VEHICLE STOLEN	04/22/2017 12:00:00 AM	21:00	...	Berkeley, CA\n(37.8690 -122.270455)
...
4945	17091287	VANDALISM	07/01/2017 12:00:00 AM	08:00	...	Berkeley, CA\n(37.8690 -122.270455)
4947	17038382	BURGLARY RESIDENTIAL	06/30/2017 12:00:00 AM	15:00	...	Berkeley, CA\n(37.8690 -122.270455)
5167	17091632	VANDALISM	08/15/2017 12:00:00 AM	23:30	...	Berkeley, CA\n(37.8690 -122.270455)

27 rows × 11 columns

It looks like 27 calls didn't have a recorded address in BLKADDR. Unfortunately, the data description isn't very clear on how the locations were recorded. We know that all of these calls were made for events in Berkeley, so we can likely assume that the addresses for these calls were originally somewhere in Berkeley.

Are there any missing values that were filled in?

From the missing value check above we can see that the Block_Location column has Berkeley, CA recorded if the location was missing.

In addition, an inspection of the `calls` table shows us that the EVENTDT column has the correct dates but records 12am for all of its times. Instead, the times are in the EVENTTM column.

```
# Show the first 7 rows of the table again for reference
calls.head(7)
```

4.1 Cleaning the Calls Dataset

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	Block_Location
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	2500 LE CONTE AVE\nBerkeley, CA\n(37.876965 -)
1	17020462	THEFT FROM PERSON	04/13/2017 12:00:00 AM	08:45	...	2200 SHATTUCK AVE\nBerkeley, CA\n(37.869363 -)
2	17050275	BURGLARY AUTO	08/24/2017 12:00:00 AM	18:30	...	200 UNIVERSITY AVE\nBerkeley, CA\n(37.865491 -)
3	17019145	GUN/WEAPON	04/06/2017 12:00:00 AM	17:30	...	1900 SEVENTH ST\nBerkeley, CA\n(37.869318 -12...)
4	17044993	VEHICLE STOLEN	08/01/2017 12:00:00 AM	18:00	...	100 PARKSIDE DR\nBerkeley, CA\n(37.854247 -12...)
5	17037319	BURGLARY RESIDENTIAL	06/28/2017 12:00:00 AM	12:00	...	1500 PRINCE ST\nBerkeley, CA\n(37.851503 -122...)
6	17030791	BURGLARY RESIDENTIAL	05/30/2017 12:00:00 AM	08:45	...	300 MENLO PL\nBerkeley, CA\n(37.851503 -122...)

7 rows × 11 columns

As a data cleaning step, we want to merge the EVENTDT and EVENTTM columns to record both date and time in one field. If we define a function that takes in a DF and returns a new DF, we can later use `pd.pipe` to apply all transformations in one go.

```

def combine_event_datetimes(calls):
    combined = pd.to_datetime(
        # Combine date and time strings
        calls['EVENTDT'].str[:10] + ' ' + calls['EVENTTM'],
        infer_datetime_format=True,
    )
    return calls.assign(EVENTDTTM=combined)

# To peek at the result without mutating the calls DF:
calls.pipe(combine_event_datetimes).head(2)

```

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	BLKADDR	
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	2500 LE CONTE AVE	Berl
1	17020462	THEFT FROM PERSON	04/13/2017 12:00:00 AM	08:45	...	2200 SHATTUCK AVE	Berl

2 rows × 12 columns

Which parts of the data were entered by a human?¶

It looks like most of the data columns are machine-recorded, including the date, time, day of week, and location of the event.

In addition, the OFFENSE and CVLEGEND columns appear to contain consistent values. We can check the unique values in each column to see if anything was misspelled:

```
calls['OFFENSE'].unique()
```

```
array(['BURGLARY AUTO', 'THEFT FROM PERSON', 'GUN/WEAPON',
       'VEHICLE STOLEN', 'BURGLARY RESIDENTIAL', 'VANDALISM',
       'DISTURBANCE', 'THEFT MISD. (UNDER $950)', 'THEFT FROM
AUTO',
       'DOMESTIC VIOLENCE', 'THEFT FELONY (OVER $950)', 'ALCOHOL
OFFENSE',
       'MISSING JUVENILE', 'ROBBERY', 'IDENTITY THEFT',
       'ASSAULT/BATTERY MISD.', '2ND RESPONSE', 'BRANDISHING',
       'MISSING ADULT', 'NARCOTICS', 'FRAUD/FORGERY',
       'ASSAULT/BATTERY FEL.', 'BURGLARY COMMERCIAL', 'MUNICIPAL
CODE',
       'ARSON', 'SEXUAL ASSAULT FEL.', 'VEHICLE RECOVERED',
       'SEXUAL ASSAULT MISD.', 'KIDNAPPING', 'VICE',
       'HOMICIDE'], dtype=object)
```

```
calls['CVLEGEND'].unique()
```

```
array(['BURGLARY - VEHICLE', 'LARCENY', 'WEAPONS OFFENSE',
       'MOTOR VEHICLE THEFT', 'BURGLARY - RESIDENTIAL',
       'VANDALISM',
       'DISORDERLY CONDUCT', 'LARCENY - FROM VEHICLE', 'FAMILY
OFFENSE',
       'LIQUOR LAW VIOLATION', 'MISSING PERSON', 'ROBBERY',
       'FRAUD',
       'ASSAULT', 'NOISE VIOLATION', 'DRUG VIOLATION',
       'BURGLARY - COMMERCIAL', 'ALL OTHER OFFENSES', 'ARSON',
       'SEX CRIME',
       'RECOVERED VEHICLE', 'KIDNAPPING', 'HOMICIDE'],
       dtype=object)
```

Since each value in these columns appears to be spelled correctly, we won't have to perform any corrections on these columns.

We also check the BLKADDR column for inconsistencies and find that sometimes an address is recorded (e.g. 2500 LE CONTE AVE) but other times a cross street is recorded (e.g. ALLSTON WAY & FIFTH ST). This suggests that a human entered this data in and this column will be difficult to use for analysis. Fortunately we can use the latitude and longitude of the event instead of the street address.

```
calls['BLKADDR'][[0, 5001]]
```

```
0          2500 LE CONTE AVE
5001      ALLSTON WAY & FIFTH ST
Name: BLKADDR, dtype: object
```

Final Touchups

This dataset seems almost ready for analysis. The Block_Location column seems to contain strings that record address, latitude, and longitude. We will want to separate the latitude and longitude for easier use.

```
def split_lat_lon(calls):
    return calls.join(
        calls['Block_Location']
        # Get coords from string
        .str.split('\n').str[2]
        # Remove parens from coords
        .str[1:-1]
        # Split latitude and longitude
        .str.split(', ', expand=True)
        .rename(columns={0: 'Latitude', 1: 'Longitude'})
    )

calls.pipe(split_lat_lon).head(2)
```

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	City	State
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	Berkeley	CA
1	17020462	THEFT FROM PERSON	04/13/2017 12:00:00 AM	08:45	...	Berkeley	CA

2 rows × 13 columns

Then, we can match the day of week number with its weekday:

4.1 Cleaning the Calls Dataset

```
# This DF contains the day for each number in CVDOW
day_of_week = pd.read_csv('data/cvdow.csv')
day_of_week
```

	CVDOW	Day
0	0	Sunday
1	1	Monday
2	2	Tuesday
3	3	Wednesday
4	4	Thursday
5	5	Friday
6	6	Saturday

```
def match_weekday(calls):
    return calls.merge(day_of_week, on='CVDOW')
calls.pipe(match_weekday).head(2)
```

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	BLKADDR	
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	2500 LE CONTE AVE	Ber
1	17038302	BURGLARY AUTO	07/02/2017 12:00:00 AM	22:00	...	BOWDITCH STREET & CHANNING WAY	Ber

2 rows × 12 columns

We'll drop columns we no longer need:

```
def drop_unneeded_cols(calls):
    return calls.drop(columns=['CVDOW', 'InDbDate',
    'Block_Location', 'City',
    'State', 'EVENTDT', 'EVENTTM'])
```

Finally, we'll pipe the `calls` DF through all the functions we've defined:

```
calls_final = (calls.pipe(combine_event_datetimes)
               .pipe(split_lat_lon)
               .pipe(match_weekday)
               .pipe(drop_unneeded_cols))
df_interact(calls_final)
```

Show Widget

(5508 rows, 8 columns) total

The Calls dataset is now ready for further data analysis. In the next section, we will clean the Stops dataset.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Cleaning The Stops Dataset](#)
 - [Are there missing values?](#)
 - [Are there any missing values that were filled in?](#)
 - [Which parts of the data were entered by a human?](#)
- [Conclusion](#)

Cleaning The Stops Dataset

The Stops dataset ([webpage](#)) records police stops of pedestrians and vehicles. Let's prepare it for further analysis.

We can use the `head` command to display the first few lines of the file.

```
!head data/stops.json
```

```
{  
  "meta" : {  
    "view" : {  
      "id" : "6e9j-pj9p",  
      "name" : "Berkeley PD - Stop Data",  
      "attribution" : "Berkeley Police Department",  
      "averageRating" : 0,  
      "category" : "Public Safety",  
      "createdAt" : 1444171604,  
      "description" : "This data was extracted from the  
Department's Public Safety Server and covers the data beginning  
January 26, 2015. On January 26, 2015 the department began  
collecting data pursuant to General Order B-4 (issued December  
31, 2014). Under that order, officers were required to provide  
certain data after making all vehicle detentions (including  
bicycles) and pedestrian detentions (up to five persons). This  
data set lists stops by police in the categories of traffic,  
suspicious vehicle, pedestrian and bicycle stops. Incident  
number, date and time, location and disposition codes are also  
listed in this data.\r\n\r\nAddress data has been changed from a  
specific address, where applicable, and listed as the block  
where the incident occurred. Disposition codes were entered by  
officers who made the stop. These codes included the person(s)  
race, gender, age (range), reason for the stop, enforcement  
action taken, and whether or not a search was  
conducted.\r\n\r\nThe officers of the Berkeley Police Department  
are prohibited from biased based policing, which is defined as  
any police-initiated action that relies on the race, ethnicity,  
or national origin rather than the behavior of an individual or  
information that leads the police to a particular individual who  
has been identified as being engaged in criminal activity.",
```

The `stops.json` file is clearly not a CSV file. In this case, the file contains data in the JSON (JavaScript Object Notation) format, a commonly used data format where data is recorded in a dictionary format. Python's [json module](#) makes reading in this file as a dictionary simple.

```
import json

# Note that this could cause our computer to run out of memory
# if the file
# is large. In this case, we've verified that the file is small
# enough to
# read in beforehand.
with open('data/stops.json') as f:
    stops_dict = json.load(f)

stops_dict.keys()
```

```
dict_keys(['meta', 'data'])
```

Note that `stops_dict` is a Python dictionary, so displaying it will display the entire dataset in the notebook. This could cause the browser to crash, so we only display the keys of the dictionary above. To peek at the data without potentially crashing the browser, we can print the dictionary to a string and only output some of the first characters of the string.

```
from pprint import pformat

def print_dict(dictionary, num_chars=1000):
    print(pformat(dictionary)[:num_chars])

print_dict(stops_dict['meta'])
```

```
{'view': {'attribution': 'Berkeley Police Department',
    'averageRating': 0,
    'category': 'Public Safety',
    'columns': [{dataTypeName': 'meta_data',
        'fieldName': ':sid',
        'flags': ['hidden'],
        'format': {},
        'id': -1,
        'name': 'sid',
        'position': 0,
        'renderTypeName': 'meta_data'},
    {'dataTypeName': 'meta_data',
        'fieldName': ':id',
        'flags': ['hidden'],
        'format': {},
        'id': -1,
        'name': 'id',
        'position': 0,
        'renderTypeName': 'meta_data'},
    {'dataTypeName': 'meta_data',
        'fieldName': ':position',
        'flags': ['hidden'],
        'format': {}}],
```

```
print_dict(stops_dict['data'], num_chars=300)
```

```
[[1,
  '29A1B912-A0A9-4431-ADC9-FB375809C32E',
  1,
  1444146408,
  '932858',
  1444146408,
  '932858',
  None,
  '2015-00004825',
  '2015-01-26T00:10:00',
  'SAN PABLO AVE / MARIN AVE',
  'T',
  'M',
  None,
  None],
 [2,
  '1644D161-1113-4C4F-BB2E-BF780E7AE73E',
  2,
  1444146408,
  '932858',
  14]
```

We can likely deduce that the `'meta'` key in the dictionary contains a description of the data and its columns and the `'data'` contains a list of data rows. We can use this information to initialize a DataFrame.

```
# Load the data from JSON and assign column titles
stops = pd.DataFrame(
    stops_dict['data'],
    columns=[c['name'] for c in stops_dict['meta']['view']
    ['columns']])
```

stops

4.2 Cleaning the Stops Dataset

	sid	id	position	created_at	...	Incident Type	D
0	1	29A1B912-A0A9-4431-ADC9-FB375809C32E	1	1444146408	...	T	N
1	2	1644D161-1113-4C4F-BB2E-BF780E7AE73E	2	1444146408	...	T	N
2	3	5338ABAB-1C96-488D-B55F-6A47AC505872	3	1444146408	...	T	N
...
29205	31079	C2B606ED-7872-4B0B-BC9B-4EF45149F34B	31079	1496269085	...	T	B
29206	31080	8FADF18D-7FE9-441D-8709-7BFEABDACA7A	31080	1496269085	...	T	H
29207	31081	F60BD2A4-8C47-4BE7-B1C6-4934BE9DF838	31081	1496269085	...	1194	A

29208 rows × 15 columns

```
# Prints column names
stops.columns
```

```
Index(['sid', 'id', 'position', 'created_at', 'created_meta',
       'updated_at',
       'updated_meta', 'meta', 'Incident Number', 'Call
Date/Time', 'Location',
       'Incident Type', 'Dispositions', 'Location - Latitude',
       'Location - Longitude'],
      dtype='object')
```

The website contains documentation about the following columns:

Column	Description	Type
Incident Number	Number of incident created by Computer Aided Dispatch (CAD) program	Plain Text
Call Date/Time	Date and time of the incident/stop	Date & Time
Location	General location of the incident/stop	Plain Text
Incident Type	This is the occurred incident type created in the CAD program. A code signifies a traffic stop (T), suspicious vehicle stop (1196), pedestrian stop (1194) and bicycle stop (1194B).	Plain Text
Dispositions	Ordered in the following sequence: 1st Character = Race, as follows: A (Asian) B (Black) H (Hispanic) O (Other) W (White) 2nd Character = Gender, as follows: F (Female) M (Male) 3rd Character = Age Range, as follows: 1 (Less than 18) 2 (18-29) 3 (30-39), 4 (Greater than 40) 4th Character = Reason, as follows: I (Investigation) T (Traffic) R (Reasonable Suspicion) K (Probation/Parole) W (Wanted) 5th Character = Enforcement, as follows: A (Arrest) C (Citation) O (Other) W (Warning) 6th Character = Car Search, as follows: S (Search) N (No Search) Additional dispositions may also appear. They are: P - Primary case report M - MDT narrative only AR - Arrest report only (no case report submitted) IN - Incident report FC - Field Card CO - Collision investigation report MH - Emergency Psychiatric Evaluation TOW - Impounded vehicle 0 or 00000 – Officer made a stop of more than five persons	Plain Text
Location - Latitude	General latitude of the call. This data is only uploaded after January 2017	Number
Location - Longitude	General longitude of the call. This data is only uploaded after January 2017.	Number

Notice that the website doesn't contain descriptions for the first 8 columns of the `stops` table. Since these columns appear to contain metadata that we're not interested in analyzing this time, we drop them from the table.

```

columns_to_drop = ['sid', 'id', 'position', 'created_at',
'created_meta',
'updated_at', 'updated_meta', 'meta']

# This function takes in a DF and returns a DF so we can use it
for .pipe
def drop_unneeded_cols(stops):
    return stops.drop(columns=columns_to_drop)

stops.pipe(drop_unneeded_cols)

```

	Incident Number	Call Date/Time	Location	Incident Type	Dispositions	
0	2015-00004825	2015-01-26T00:10:00	SAN PABLO AVE / MARIN AVE	T	M	Nc
1	2015-00004829	2015-01-26T00:50:00	SAN PABLO AVE / CHANNING WAY	T	M	Nc
2	2015-00004831	2015-01-26T01:03:00	UNIVERSITY AVE / NINTH ST	T	M	Nc
...
29205	2017-00024245	2017-04-30T22:59:26	UNIVERSITY AVE/6TH ST	T	BM2TWN;	Nc
29206	2017-00024250	2017-04-30T23:19:27	UNIVERSITY AVE / WEST ST	T	HM4TCS;	37
29207	2017-00024254	2017-04-30T23:38:34	CHANNING WAY / BOWDITCH ST	1194	AR;	37

29208 rows × 7 columns

As with the Calls dataset, we will answer the following three questions about the Stops dataset:

1. Are there missing values in the dataset?
2. Are there any missing values that were filled in (e.g. a 999 for unknown age or 12:00am for unknown date)?

3. Which parts of the data were entered by a human?

Are there missing values?

We can clearly see that there are many missing latitude and longitudes. The data description states that these two columns are only filled in after Jan 2017.

```
# True if row contains at least one null value
null_rows = stops.isnull().any(axis=1)

stops=null_rows]
```

	Incident Number	Call Date/Time	Location	Incident Type	Dispositions	Lo La
0	2015-00004825	2015-01-26T00:10:00	SAN PABLO AVE / MARIN AVE	T	M	Nc
1	2015-00004829	2015-01-26T00:50:00	SAN PABLO AVE / CHANNING WAY	T	M	Nc
2	2015-00004831	2015-01-26T01:03:00	UNIVERSITY AVE / NINTH ST	T	M	Nc
...
29078	2017-00023764	2017-04-29T01:59:36	2180 M L KING JR WAY	1194	BM4IWN;	Nc
29180	2017-00024132	2017-04-30T12:54:23	6TH/UNI	1194	M;	Nc
29205	2017-00024245	2017-04-30T22:59:26	UNIVERSITY AVE/6TH ST	T	BM2TWN;	Nc

25067 rows × 7 columns

We can check the other columns for missing values:

```
# True if row contains at least one null value without checking
# the latitude and longitude columns
null_rows = stops.iloc[:, :-2].isnull().any(axis=1)

df_interact(stops=null_rows)
```

Show Widget

(63 rows, 7 columns) total

By browsing through the table above, we can see that all other missing values are in the Dispositions column. Unfortunately, we do not know from the data description why these Dispositions might be missing. Since only there are only 63 missing values compared to 25,000 rows in the original table, we can proceed with analysis while being mindful that these missing values could impact results.

Are there any missing values that were filled in?

It doesn't seem like any previously missing values were filled in for us. Unlike in the Calls dataset where the date and time were in separate columns, the Call Date/Time column in the Stops dataset contains both date and time.

Which parts of the data were entered by a human?

As with the Calls dataset, it looks like most of the columns in this dataset were recorded by a machine or were a category selected by a human (e.g. Incident Type).

However, the Location column doesn't have consistently entered values. Sure enough, we spot some typos in the data:

```
stops['Location'].value_counts()
```

```

2200 BLOCK SHATTUCK AVE           229
37.8693028530001~-122.272234021 213
UNIVERSITY AVE / SAN PABLO AVE    202
...
VALLEY ST / DWIGHT WAY           1
COLLEGE AVE / SIXTY-THIRD ST     1
GRIZZLY PEAK BLVD / MARIN AVE    1
Name: Location, Length: 6393, dtype: int64

```

What a mess! It looks like sometimes an address was entered, sometimes a cross-street, and other times a latitude-longitude pair. Unfortunately, we don't have very complete latitude-longitude data to use in place of this column. We may have to manually clean this column if we want to use locations for future analysis.

We can also check the Dispositions column:

```

dispositions = stops['Dispositions'].value_counts()

# Outputs a slider to pan through the unique Dispositions in
# order of how often they appear
interact(lambda row=0: dispositions.iloc[row:row+7],
          row=(0, len(dispositions), 7))

```

Show Widget

```
<function __main__.<lambda>>
```

The Dispositions columns also contains inconsistencies. For example, some dispositions start with a space, some end with a semicolon, and some contain multiple entries. The variety of values suggests that this field contains human-entered values and should be treated with caution.

```

# Strange values...
dispositions.iloc[[0, 20, 30, 266, 1027]]

```

```
M           1683
M;          238
M           176
HF4TWN;     14
OM4KWS      1
Name: Dispositions, dtype: int64
```

In addition, the most common disposition is `M` which isn't a permitted first character in the `Dispositions` column. This could mean that the format of the column changed over time or that officers are allowed to enter in the disposition without matching the format in the data description. In any case, the column will be challenging to work with.

We can take some simple steps to clean the `Dispositions` column by removing leading and trailing whitespace, removing trailing semi-colons, and replacing the remaining semi-colons with commas.

```
def clean_dispositions(stops):
    cleaned = (stops['Dispositions']
               .str.strip()
               .str.rstrip(';')
               .str.replace(';', ','))
    return stops.assign(Dispositions=cleaned)
```

As before, we can now pipe the `stops` DF through the cleaning functions we've defined:

```
stops_final = (stops
               .pipe(drop_unneeded_cols)
               .pipe(clean_dispositions))
df_interact(stops_final)
```

Show Widget

```
(29208 rows, 7 columns) total
```

Conclusion

As these two datasets have shown, data cleaning can often be both difficult and tedious. Cleaning 100% of the data often takes too long, but not cleaning the data at all results in faulty conclusions; we have to weigh our options and strike a balance each time we encounter a new dataset.

The decisions made during data cleaning impact all future analyses. For example, we chose not to clean the Location column of the Stops dataset so we should treat that column with caution. Each decision made during data cleaning should be carefully documented for future reference, preferably in a notebook so that both code and explanations appear together.

Exploratory Data Analysis

Exploratory data analysis is an attitude, a state of flexibility, a willingness to look for those things that we believe are not there, as well as those we believe to be there.

— [John Tukey](#)

In Exploratory Data Analysis (EDA), the third step of the data science lifecycle, we summarize, visualize, and transform the data in order to understand it more deeply. In particular, through EDA we identify potential issues in the data and discover trends that inform further analyses.

We seek to understand the following properties about our data:

1. Structure: the format of our data file.
2. Granularity: how fine or coarse each row and column is.
3. Scope: how (in)complete our data are.
4. Temporality: how the data are situated in time.
5. Faithfulness: how well the data captures "reality".

Although we introduce data cleaning and EDA separately to help organize this book, in practice you will often switch between the two. For example, visualizing a column may show misformatted values that you should use data cleaning techniques to process. With this in mind, we return to the Berkeley Police Department datasets for exploration.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Structure](#)
- [Joins](#)
- [Structure Checklist](#)

Structure

The structure of a dataset refers to the "shape" of the data files. At a basic level, this refers to the format that the data are entered in. For example, we saw that the Calls dataset is a comma-separated values file:

```
!head data/Berkeley_PD_-_Calls_for_Service.csv
```

```
CASENO,OFFENSE,EVENTDT,EVENTTM,CVLEGEND,CVDOW,InDbDate,Block_Location,BLKADDR,City,State  
17091420,BURGLARY AUTO,07/23/2017 12:00:00 AM,06:00,BURGLARY -  
VEHICLE,0,08/29/2017 08:28:05 AM,"2500 LE CONTE AVE  
Berkeley, CA  
(37.876965, -122.260544)",2500 LE CONTE AVE,Berkeley,CA  
17020462,THEFT FROM PERSON,04/13/2017 12:00:00  
AM,08:45,LARCENY,4,08/29/2017 08:28:00 AM,"2200 SHATTUCK AVE  
Berkeley, CA  
(37.869363, -122.268028)",2200 SHATTUCK AVE,Berkeley,CA  
17050275,BURGLARY AUTO,08/24/2017 12:00:00 AM,18:30,BURGLARY -  
VEHICLE,4,08/29/2017 08:28:06 AM,"200 UNIVERSITY AVE  
Berkeley, CA  
(37.865491, -122.310065)",200 UNIVERSITY AVE,Berkeley,CA
```

The Stops dataset, on the other hand, is a JSON (JavaScript Object Notation) file.

```
# Show first and last 5 lines of file
!head -n 5 data/stops.json
!echo '...'
!tail -n 5 data/stops.json
```

```
{
  "meta" : {
    "view" : {
      "id" : "6e9j-pj9p",
      "name" : "Berkeley PD - Stop Data",
      ...
      , [ 31079, "C2B606ED-7872-4B0B-BC9B-4EF45149F34B", 31079,
        1496269085, "932858", 1496269085, "932858", null, "2017-
        00024245", "2017-04-30T22:59:26", " UNIVERSITY AVE/6TH ST", "T",
        "BM2TWN; ", null, null ]
      , [ 31080, "8FADF18D-7FE9-441D-8709-7BFEABDACA7A", 31080,
        1496269085, "932858", 1496269085, "932858", null, "2017-
        00024250", "2017-04-30T23:19:27", " UNIVERSITY AVE / WEST ST",
        "T", "HM4TCS; ", "37.8698757000001", "-122.286550846" ]
      , [ 31081, "F60BD2A4-8C47-4BE7-B1C6-4934BE9DF838", 31081,
        1496269085, "932858", 1496269085, "932858", null, "2017-
        00024254", "2017-04-30T23:38:34", " CHANNING WAY / BOWDITCH
        ST", "1194", "AR; ", "37.867207539", "-122.256529377" ]
      ]
    }
  }
```

Of course, there are many other types of data formats. Here is a list of the most common formats:

- Comma-Separated Values (CSV) and Tab-Separated Values (TSV). These files contain tabular data delimited by either a comma for CSV or a tab character (\t) for TSV. These files are typically easy to work with because the data are entered in a similar format to DataFrames.
- JavaScript Object Notation (JSON). These files contain data in a nested dictionary format. Typically we have to read in the entire file as a Python dict and then figure out how to extract fields for a DataFrame from the dict.
- eXtensible Markup Language (XML) or HyperText Markup Language (HTML). These files also contain data in a nested format, for example:

```
<?xml version="1.0" encoding="UTF-8"?>
<note>
  <to>Tove</to>
  <from>Jani</from>
  <heading>Reminder</heading>
  <body>Don't forget me this weekend!</body>
</note>
```

In a later chapter we will use XPath to extract data from these types of files.

- Log data. Many applications will output some data as they run in an unstructured text format, for example:

```
2005-03-23 23:47:11,663 - sa - INFO - creating an instance of aux_module.Aux
2005-03-23 23:47:11,665 - sa.aux.Aux - INFO - creating an instance of Aux
2005-03-23 23:47:11,665 - sa - INFO - created an instance of aux_module.Aux
2005-03-23 23:47:11,668 - sa - INFO - calling aux_module.Aux.do_something
2005-03-23 23:47:11,668 - sa.aux.Aux - INFO - doing something
```

In a later chapter we will use Regular Expressions to extract data from these types of files.

Joins

Data will often be split across multiple tables. For example, one table can describe some people's personal information while another will contain their emails:

```
people = pd.DataFrame(
    [[{"Name": "Joey", "Color": "blue", "Number": 42, "Sex": "M"}, {"Name": "Weiwei", "Color": "blue", "Number": 50, "Sex": "F"}, {"Name": "Joey", "Color": "green", "Number": 8, "Sex": "M"}, {"Name": "Karina", "Color": "green", "Number": 7, "Sex": "F"}, {"Name": "Nhi", "Color": "blue", "Number": 3, "Sex": "F"}, {"Name": "Sam", "Color": "pink", "Number": -42, "Sex": "M"}],
    columns = ["Name", "Color", "Number", "Sex"])
```

people

	Name	Color	Number	Sex
0	Joey	blue	42	M
1	Weiwei	blue	50	F
2	Joey	green	8	M
3	Karina	green	7	F
4	Fernando	pink	-9	M
5	Nhi	blue	3	F
6	Sam	pink	-42	M

```
email = pd.DataFrame(
    [[ "Deb", "deborah_nolan@berkeley.edu"],
     [ "Sam", "samlau95@berkeley.edu"],
     [ "John", "doe@nope.com"],
     [ "Joey", "jegonzal@cs.berkeley.edu"],
     [ "Weiwei", "weiwzhang@berkeley.edu"],
     [ "Weiwei", "weiwzhang+123@berkeley.edu"],
     [ "Karina", "kgoot@berkeley.edu"]],
    columns = ["User Name", "Email"])
```

```
email
```

	User Name	Email
0	Deb	deborah_nolan@berkeley.edu
1	Sam	samlau95@berkeley.edu
2	John	doe@nope.com
3	Joey	jegonzal@cs.berkeley.edu
4	Weiwei	weiwzhang@berkeley.edu
5	Weiwei	weiwzhang+123@berkeley.edu
6	Karina	kgoot@berkeley.edu

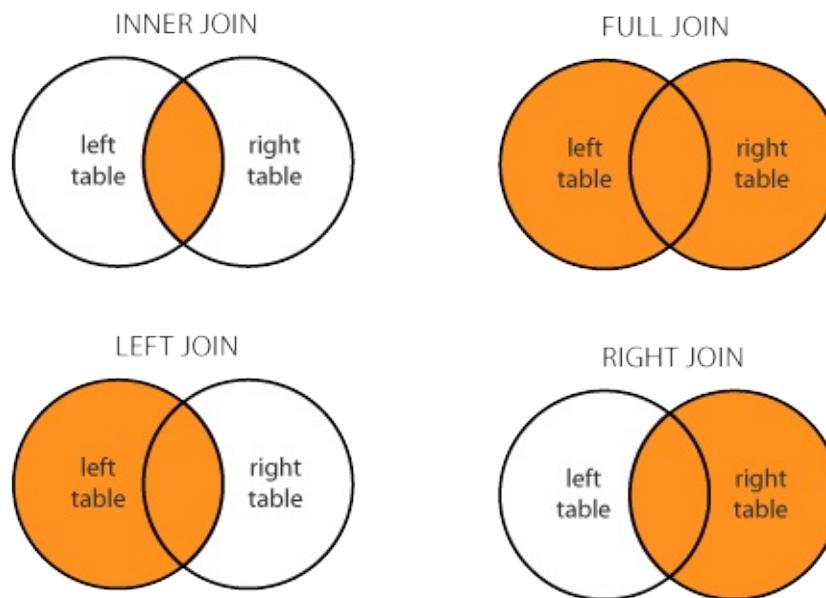
To match up each person with his or her email, we can join the two tables on the columns that contain the usernames. We must then decide what to do about people that appear in one table but not the other. For example, Fernando appears in the `people` table but not the

`email` table. We have several types of joins for each strategy of matching missing values. One of the more common joins is the *inner join*, where any row that doesn't have a match is dropped in the final result:

```
# Fernando, Nhi, Deb, and John don't appear
people.merge(email, how='inner', left_on='Name', right_on='User
Name')
```

	Name	Color	Number	Sex	User Name	Email
0	Joey	blue	42	M	Joey	jegonzal@cs.berkeley.edu
1	Joey	green	8	M	Joey	jegonzal@cs.berkeley.edu
2	Weiwei	blue	50	F	Weiwei	weiwzhang@berkeley.edu
3	Weiwei	blue	50	F	Weiwei	weiwzhang+123@berkeley.edu
4	Karina	green	7	F	Karina	kgoot@berkeley.edu
5	Sam	pink	-42	M	Sam	samlau95@berkeley.edu

There are four basic joins that we use most often: inner, full (sometimes called "outer"), left, and right joins. Below is a diagram to show the difference between these types of joins.



Use the dropdown menu below to show the result of the four different types of joins on the `people` and `email` tables. Notice which rows contain NaN values for outer, left, and right joins.

Show Widget

Structure Checklist

You should have answers to the following questions after looking at the structure of your datasets. We will answer them for the Calls and Stops datasets.

Are the data in a standard format or encoding?

Standard formats include:

- Tabular data: CSV, TSV, Excel, SQL
- Nested data: JSON, XML

The Calls dataset came in the CSV format while the Stops dataset came in the JSON format.

Are the data organized in records (e.g. rows)? If not, can we define records by parsing the data?

The Calls dataset came in rows; we extracted records from the Stops dataset.

Are the data nested? If so, can we reasonably unnest the data?

The Calls dataset wasn't nested; we didn't have to work too hard to unnest data from the Stops dataset.

Do the data reference other data? If so, can we join the data?

The Calls dataset references the day of week table. Joining those two tables gives us the day of week for each incident in the dataset. The Stops dataset had no obvious references.

What are the fields (e.g. columns) in each record? What is the type of each column?

The fields for the Calls and Stops datasets are described in the Data Cleaning sections for each dataset.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Granularity](#)
- [Granularity Checklist](#)

Granularity

The granularity of your data is what each record in your data represents. For example, in the Calls dataset each record represents a single case of a police call.

	CASENO	OFFENSE	CVLEGEND	BLKADDR	EVENTDTTM
0	17091420	BURGLARY AUTO	BURGLARY - VEHICLE	2500 LE CONTE AVE	2017-07-23 06:00:00
1	17038302	BURGLARY AUTO	BURGLARY - VEHICLE	BOWDITCH STREET & CHANNING WAY	2017-07-02 22:00:00
2	17049346	THEFT MISD. (UNDER \$950)	LARCENY	2900 CHANNING WAY	2017-08-20 23:20:00
3	17091319	THEFT MISD. (UNDER \$950)	LARCENY	2100 RUSSELL ST	2017-07-09 04:15:00
4	17044238	DISTURBANCE	DISORDERLY CONDUCT	TELEGRAPH AVENUE & DURANT AVE	2017-07-30 01:16:00

In the Stops dataset, each record represents a single incident of a police stop.

	Incident Number	Call Date/Time	Location	Incident Type	Dispositions	Location - Latitude
0	2015-00004825	2015-01-26 00:10:00	SAN PABLO AVE / MARIN AVE	T	M	NaN
1	2015-00004829	2015-01-26 00:50:00	SAN PABLO AVE / CHANNING WAY	T	M	NaN
2	2015-00004831	2015-01-26 01:03:00	UNIVERSITY AVE / NINTH ST	T	M	NaN
3	2015-00004848	2015-01-26 07:16:00	2000 BLOCK BERKELEY WAY	1194	BM4ICN	NaN
4	2015-00004849	2015-01-26 07:43:00	1700 BLOCK SAN PABLO AVE	1194	BM4ICN	NaN

On the other hand, we could have received the Stops data in the following format:

	Num Incidents
Call Date/Time	
2015-01-26	46
2015-01-27	57
2015-01-28	56
...	...
2017-04-28	82
2017-04-29	86
2017-04-30	59

825 rows × 1 columns

In this case, each record in the table corresponds to a single date instead of a single incident. We would describe this table as having a coarser granularity than the one above. It's important to know the granularity of your data because it determines what kind of analyses you can perform. Generally speaking, too fine of a granularity is better than too coarse; while we can use grouping and pivoting to change a fine granularity to a coarse one, we have few tools to go from coarse to fine.

Granularity Checklist

You should have answers to the following questions after looking at the granularity of your datasets. We will answer them for the Calls and Stops datasets.

What does a record represent?

In the Calls dataset, each record represents a single case of a police call. In the Stops dataset, each record represents a single incident of a police stop.

Do all records capture granularity at the same level? (Sometimes a table will contain summary rows.)

Yes, for both Calls and Stops datasets.

If the data were aggregated, how was the aggregation performed? Sampling and averaging are common aggregations.

No aggregations were performed as far as we can tell for the datasets. We do keep in mind that in both datasets, the location is entered as a block location instead of a specific address.

What kinds of aggregations can we perform on the data?

For example, it's often useful to aggregate individual people to demographic groups or individual events to totals across time.

In this case, we can aggregate across various granularities of date or time. For example, we can find the most common hour of day for incidents with aggregation. We might also be able to aggregate across event locations to find the regions of Berkeley with the most incidents.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Scope](#)

Scope

The scope of the dataset refers to the coverage of the dataset in relation to what we are interested in analyzing. We seek to answer the following question about our data scope:

Does the data cover the topic of interest?

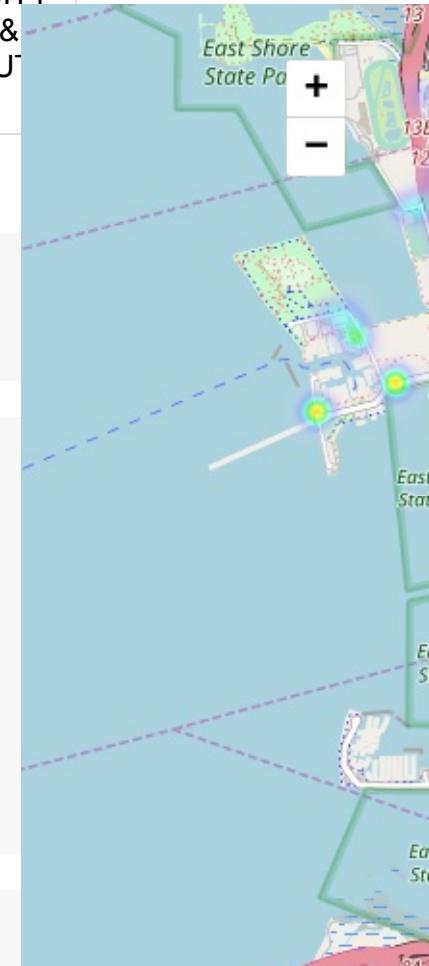
For example, the Calls and Stops datasets contain call and stop incidents made in Berkeley. If we are interested in crime incidents in the state of California, however, these datasets will be too limited in scope.

In general, larger scope is more useful than smaller scope since we can filter larger scope down to a smaller scope but often can't go from smaller scope to larger scope. For example, if we had a dataset of police stops in the United States we could subset the dataset to investigate Berkeley.

Keep in mind that scope is a broad term not always used to describe geographic location. For example, it can also refer to time coverage — the Calls dataset only contains data for a 180 day period.

We will often address the scope of the dataset during the investigation of the data generation process and confirm the dataset's scope during EDA. Let's confirm the geographic and time scope of the Calls dataset.

```
calls
```

	CASENO	OFFENSE	CVLEGEND	BLKADDR	EVENTDTT
0	17091420	BURGLARY AUTO	BURGLARY - VEHICLE	2500 LE CONTE AVE	2017-07-23 06:00:00
1	17038302	BURGLARY AUTO	BURGLARY - VEHICLE	BOWDITCH STREET & CHANNING WAY	2017-07-02 22:00:00
2	17049346	THEFT MISD. (UNDER \$950)	LARCENY	2900 CHANNING WAY	2017-08-20 23:20:00
...
5505	17021604	IDENTITY THEFT	FRAUD	100 MONTROSE RD	2017-03-31 00:00:00
5506	17033201	DISTURBANCE	DISORDERLY CONDUCT	2300 COLLEGE AVE	2017-06-09 22:34:00
5507	17047247	BURGLARY AUTO	BURGLARY - VEHICLE	UNIVERSITY AVENUE & CHESTNUT ST	

5508 rows × 8 columns

```
# Shows earliest and latest dates in calls
calls['EVENTDTTM'].dt.date.sort_values()
```

1384	2017-03-02
1264	2017-03-02
1408	2017-03-02
...	
3516	2017-08-28
3409	2017-08-28
3631	2017-08-28

Name: EVENTDTTM, Length: 5508, dtype: object

```
calls['EVENTDTTM'].dt.date.max() -
calls['EVENTDTTM'].dt.date.min()
```

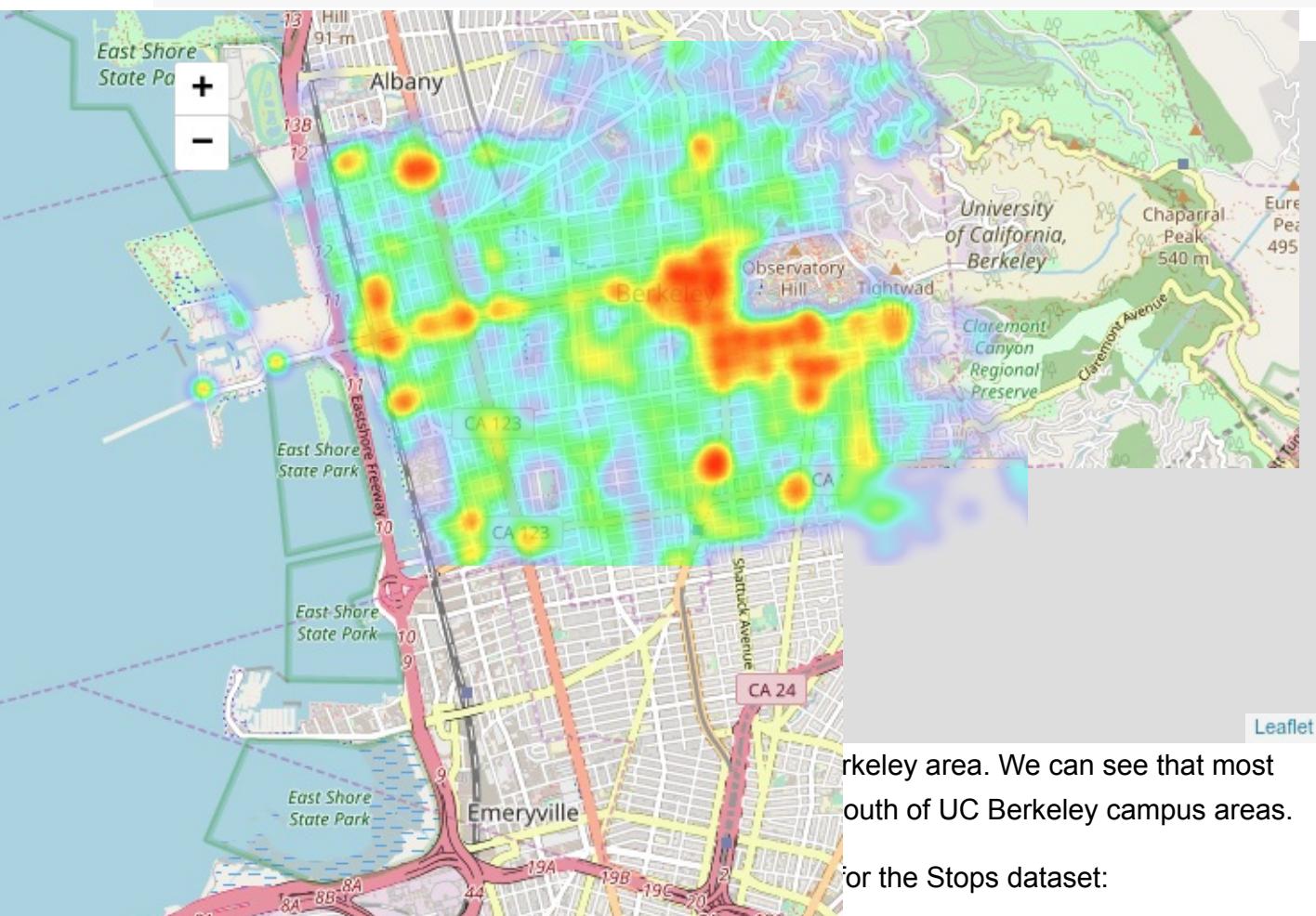
```
datetime.timedelta(179)
```

The table contains data for a time period of 179 days which is close enough to the 180 day time period in the data description that we can suppose there were no calls on either April 14st, 2017 or August 29, 2017.

To check the geographic scope, we can use a map:

```
import folium # Use the Folium Javascript Map Library
import folium.plugins

SF_COORDINATES = (37.87, -122.28)
sf_map = folium.Map(location=SF_COORDINATES, zoom_start=13)
locs = calls[['Latitude',
'Longitude']].astype('float').dropna().as_matrix()
heatmap = folium.plugins.HeatMap(locs.tolist(), radius = 10)
sf_map.add_child(heatmap)
```



keley area. We can see that most
outh of UC Berkeley campus areas.

for the Stops dataset:

```
stops
```

	Incident Number	Call Date/Time	Location	Incident Type	Dispositions	Loc - La
0	2015-00004825	2015-01-26 00:10:00	SAN PABLO AVE / MARIN AVE	T	M	Nan
1	2015-00004829	2015-01-26 00:50:00	SAN PABLO AVE / CHANNING WAY	T	M	Nan
2	2015-00004831	2015-01-26 01:03:00	UNIVERSITY AVE / NINTH ST	T	M	
...	
29205	2017-00024245	2017-04-30 22:59:26	UNIVERSITY AVE/6TH ST	T	BM2	
29206	2017-00024250	2017-04-30 23:19:27	UNIVERSITY AVE / WEST ST	T	HM4	
29207	2017-00024254	2017-04-30 23:38:34	CHANNING WAY / BOWDITCH ST	1194	AR	

29208 rows × 7 columns

```
stops['Call Date/Time'].dt.date.sort_values()
```

0	2015-01-26
25	2015-01-26
26	2015-01-26
...	
29175	2017-04-30
29177	2017-04-30
29207	2017-04-30

```
Name: Call Date/Time, Length: 29208, dtype: object
```

As promised, the data collection begins on January 26th, 2015. It looks like the data were downloaded somewhere around the beginning of May 2017 since the dates stop on April 30th, 2017. Let's draw a map to see the geographic data:

```
SF_COORDINATES = (37.87, -122.28)
sf_map = folium.Map(location=SFCORDINATES, zoom_start=13)
locs = stops[['Location - Latitude', 'Location - Longitude']].astype('float').dropna().as_matrix()
heatmap = folium.plugins.HeatMap(locs.tolist(), radius = 10)
sf_map.add_child(heatmap)
```



[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Temporality](#)

Temporality

Temporality refers to how the data are situated in time and specifically to the date and time fields in the dataset. We seek to understand the following traits about these fields:

What is the meaning of the date and time fields in the dataset?

In the Calls and Stops dataset, the datetime fields represent when the call or stop was made by the police. However, the Stops dataset also originally had a datetime field recording when the case was entered into the database which we took out during data cleaning since we didn't think it would be useful for analysis.

In addition, we should be careful to note the timezone and Daylight Savings for datetime fields especially when dealing with data that comes from multiple locations.

What representation do the date and time fields have in the data?

Although the US uses the MM/DD/YYYY format, many other countries use the DD/MM/YYYY format. There are still more formats in use around the world and it's important to recognize these differences when analyzing data.

In the Calls and Stops dataset, the dates came in the MM/DD/YYYY format.

Are there strange timestamps that might represent null values?

Some programs use placeholder datetimes instead of null values. For example, Excel's default date is Jan 1st, 1990 and on Excel for Mac, it's Jan 1st, 1904. Many applications will generate a default datetime of 12:00am Jan 1st, 1970 or 11:59pm Dec 31st, 1969 since this is the [Unix Epoch for timestamps](#). If you notice multiple instances of these timestamps in your data, you should take caution and double check your data sources. Neither Calls nor Stops dataset contain any of these suspicious values.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Faithfulness](#)

Faithfulness

We describe a dataset as "faithful" if we believe it accurately captures reality. Typically, untrustworthy datasets contain:

Unrealistic or incorrect values

For example, dates in the future, locations that don't exist, negative counts, or large outliers.

Violations of obvious dependencies

For example, age and birthday for individuals don't match.

Hand-entered data

As we have seen, these are typically filled with spelling errors and inconsistencies.

Clear signs of data falsification

For example, repeated names, fake looking email addresses, or repeated use of uncommon names or fields.

Notice the many similarities to data cleaning. As we have mentioned, we often go back and forth between data cleaning and EDA, especially when determining data faithfulness. For example, visualizations often help us identify strange entries in the data.

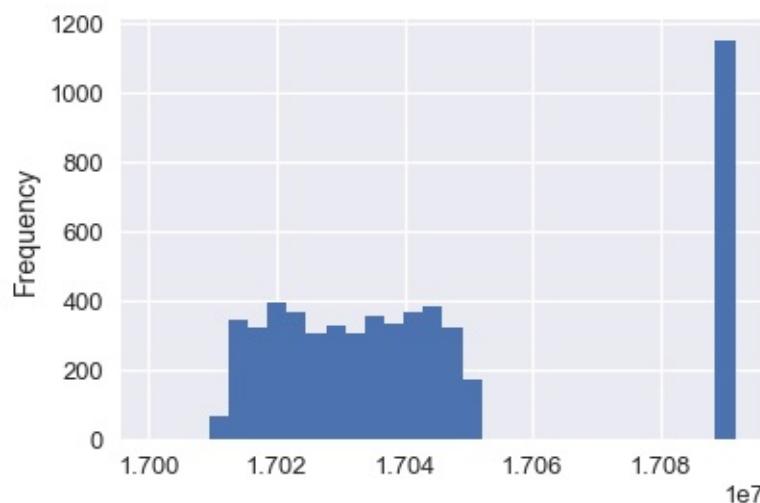
```
calls = pd.read_csv('data/calls.csv')
calls.head()
```

	CASENO	OFFENSE	EVENTDT	EVENTTM	...	BLKADDR
0	17091420	BURGLARY AUTO	07/23/2017 12:00:00 AM	06:00	...	2500 LE CONTE AVE
1	17038302	BURGLARY AUTO	07/02/2017 12:00:00 AM	22:00	...	BOWDITCH STREET & CHANNING WAY
2	17049346	THEFT MISD. (UNDER \$950)	08/20/2017 12:00:00 AM	23:20	...	2900 CHANNING WAY
3	17091319	THEFT MISD. (UNDER \$950)	07/09/2017 12:00:00 AM	04:15	...	2100 RUSSELL ST
4	17044238	DISTURBANCE	07/30/2017 12:00:00 AM	01:16	...	TELEGRAPH AVENUE & DURANT AVE

5 rows × 9 columns

```
calls['CASENO'].plot.hist(bins=30)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a1ebb2898>
```



Notice the unexpected clusters at 17030000 and 17090000. By plotting the distribution of case numbers, we can quickly see anomalies in the data. In this case, we might guess that two different teams of police use different sets of case numbers for their calls.

Exploring the data often reveals anomalies; if fixable, we can then apply data cleaning techniques.

Data Visualization

There is a magic in graphs. The profile of a curve reveals in a flash a whole situation — the life history of an epidemic, a panic, or an era of prosperity. The curve informs the mind, awakens the imagination, convinces.

— Henry D. Hubbard

Data visualization is an essential tool for data science at every step of analysis, from data cleaning to EDA to communicating conclusions and predictions. Because human minds are highly developed for visual perception, a well-chosen plot can often reveal trends and anomalies in the data much more efficiently than a textual description.

To effectively use data visualizations, you must be proficient with both the programming tools to generate plots and the principles of visualization. In this chapter we will introduce `seaborn` and `matplotlib`, our tools of choice for creating plots. We will also learn how to spot misleading visualizations and how to improve visualizations using data transformations, smoothing, and dimensionality reduction.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- Visualizing Quantitative Data
 - Histograms
 - Box plots
 - Brief Aside on Using Seaborn
 - Scatter Plots

Visualizing Quantitative Data

We generally use different types of charts to visualize quantitative (numerical) data and qualitative (ordinal or nominal) data.

For quantitative data, we most often use histograms, box plots, and scatter plots.

We can use the [seaborn plotting library](#) to create these plots in Python. We will use a dataset containing information about passengers aboard the Titanic.

```
# Import seaborn and apply its plotting styles
import seaborn as sns
sns.set()
```

```
# Load the dataset and drop N/A values to make plot function
# calls simpler
ti = sns.load_dataset('titanic').dropna().reset_index(drop=True)

# This table is too large to fit onto a page so we'll output
# sliders to
# pan through different sections.
df_interact(ti)
```

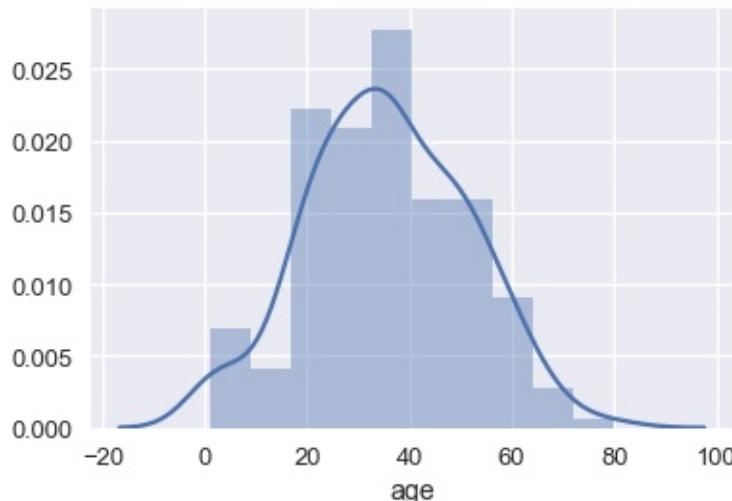
[Show Widget](#)

(182 rows, 15 columns) total

Histograms

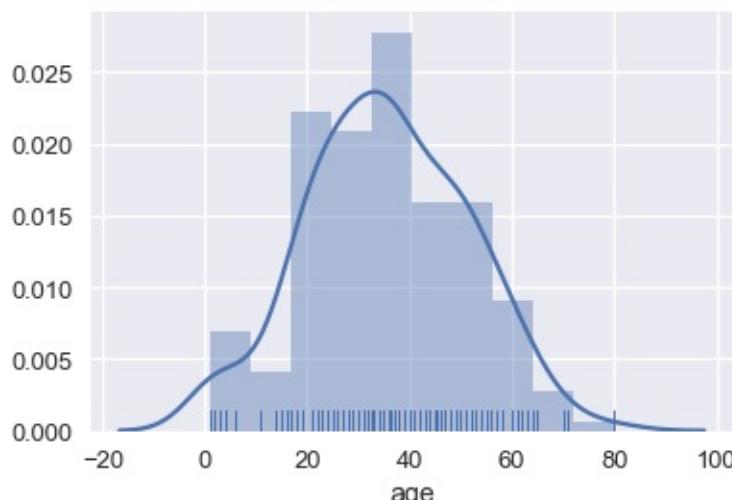
We can see that the dataset contains one row for every passenger. Each row includes the age of the passenger and the amount the passenger paid for a ticket. Let's visualize the ages using a histogram. We can use seaborn's `distplot` function:

```
# Adding a semi-colon at the end tells Jupyter not to output the
# usual <matplotlib.axes._subplots.AxesSubplot> line
sns.distplot(ti['age']);
```



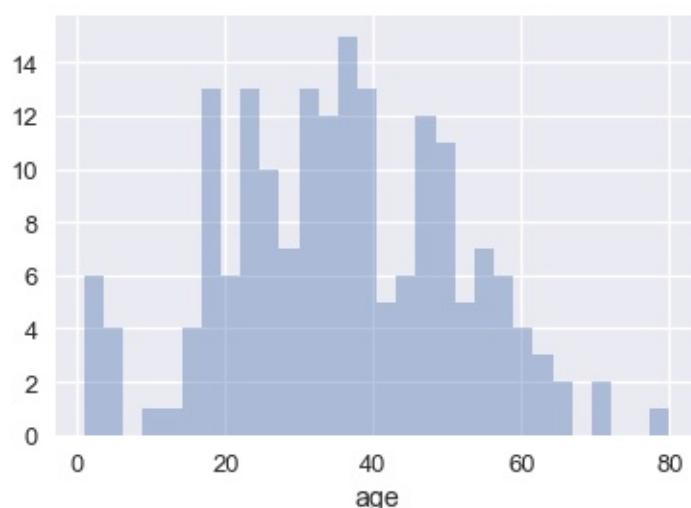
By default, seaborn's `distplot` function will output a smoothed curve that roughly fits the distribution. We can also add a rugplot which marks each individual point on the x-axis:

```
sns.distplot(ti['age'], rug=True);
```



We can also plot the distribution itself. Adjusting the number of bins shows that there were a number of children on board.

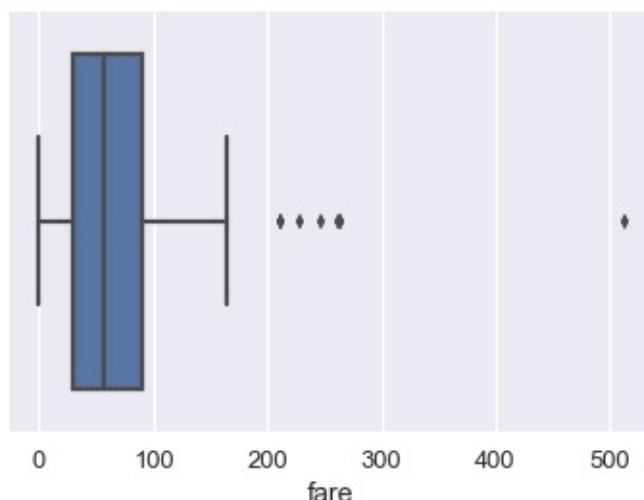
```
sns.distplot(ti['age'], kde=False, bins=30);
```



Box plots

Box plots are a convenient way to see where most of the data lie. Typically, we use the 25th and 75th percentiles of the data as the start and endpoints of the box and draw a line within the box for the 50th percentile (the median). We draw two "whiskers" that extend to show the remaining data except outliers, which are marked as individual points outside the whiskers.

```
sns.boxplot(x='fare', data=ti);
```



We typically use the Inter-Quartile Range (IQR) to determine which points are considered outliers for the box plot. The IQR is the difference between the 75th percentile of the data and the 25th percentile.

```
lower, upper = np.percentile(ti['fare'], [25, 75])
iqr = upper - lower
iqr
```

```
60.29999999999997
```

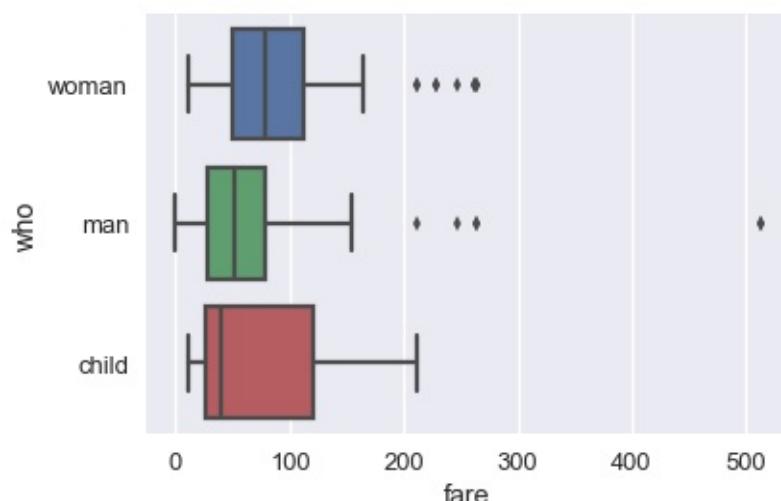
Values greater than $1.5 \times \text{IQR}$ above the 75th percentile and less than $1.5 \times \text{IQR}$ below the 25th percentile are considered outliers and we can see them marked individually on the boxplot above:

```
upper_cutoff = upper + 1.5 * iqr
lower_cutoff = lower - 1.5 * iqr
upper_cutoff, lower_cutoff
```

```
(180.449999999999, -60.7499999999986)
```

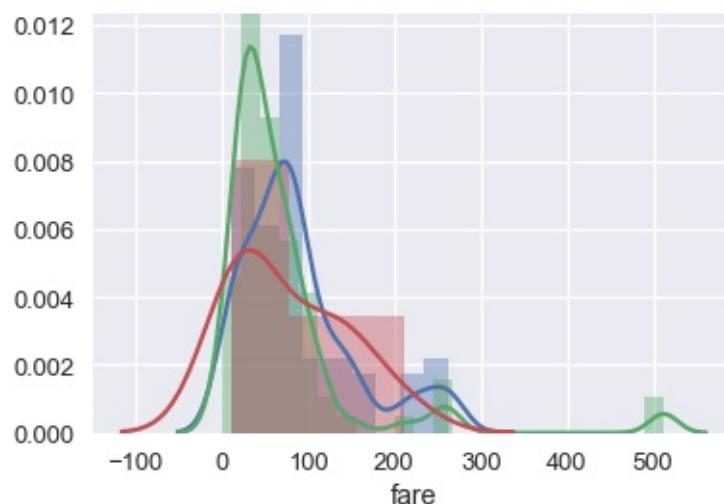
Although histograms show the entire distribution at once, box plots are often easier to understand when we split the data by different categories. For example, we can make one box plot for each passenger type:

```
sns.boxplot(x='fare', y='who', data=ti);
```



The separate box plots are much easier to understand than the overlaid histogram below which plots the same data:

```
sns.distplot(ti.loc[ti['who'] == 'woman', 'fare'])
sns.distplot(ti.loc[ti['who'] == 'man', 'fare'])
sns.distplot(ti.loc[ti['who'] == 'child', 'fare']);
```



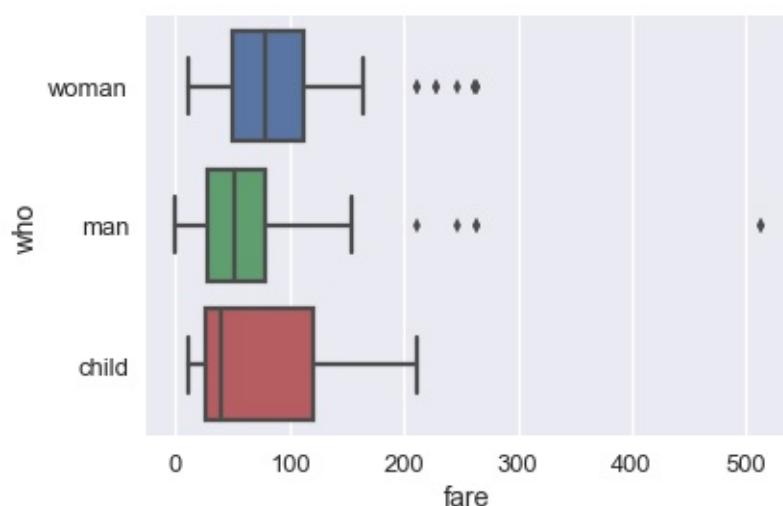
Brief Aside on Using Seaborn

You may have noticed that the `boxplot` call to make separate box plots for the `who` column was simpler than the equivalent code to make an overlaid histogram. Although `sns.distplot` takes in an array or Series of data, most other seaborn functions allow you to pass in a DataFrame and specify which column to plot on the x and y axes. For example:

```
# Plots the `fare` column of the `ti` DF on the x-axis
sns.boxplot(x='fare', data=ti);
```

When the column is categorical (the `'who'` column contained `'woman'`, `'man'`, and `'child'`), seaborn will automatically split the data by category before plotting. This means we don't have to filter out each category ourselves like we did for `sns.distplot`.

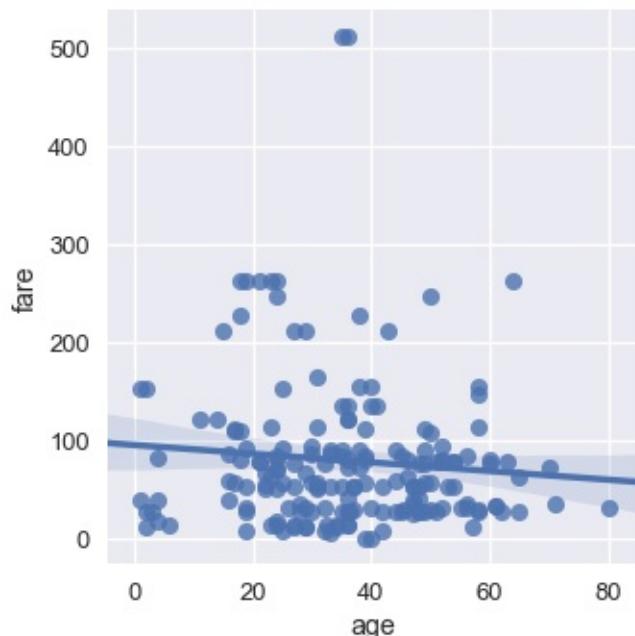
```
# fare (numerical) on the x-axis,
# who (nominal) on the y-axis
sns.boxplot(x='fare', y='who', data=ti);
```



Scatter Plots

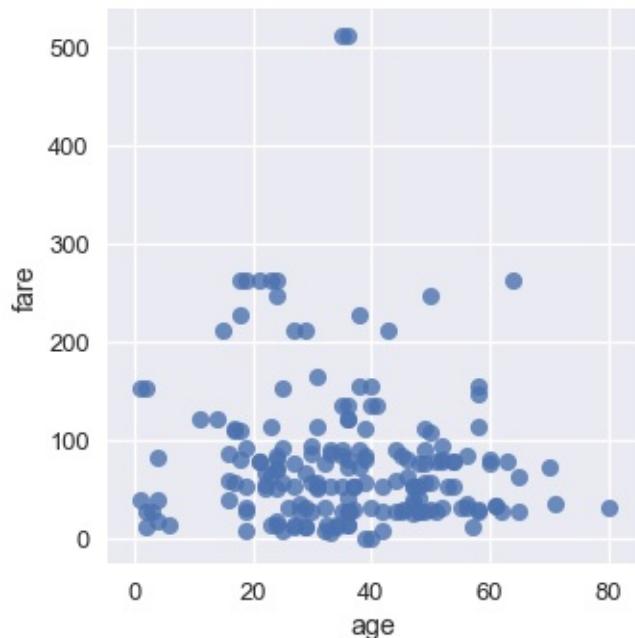
Scatter plots are used to compare two quantitative variables. We can compare the `age` and `fare` columns of our Titanic dataset using a scatter plot.

```
sns.lmplot(x='age', y='fare', data=ti);
```



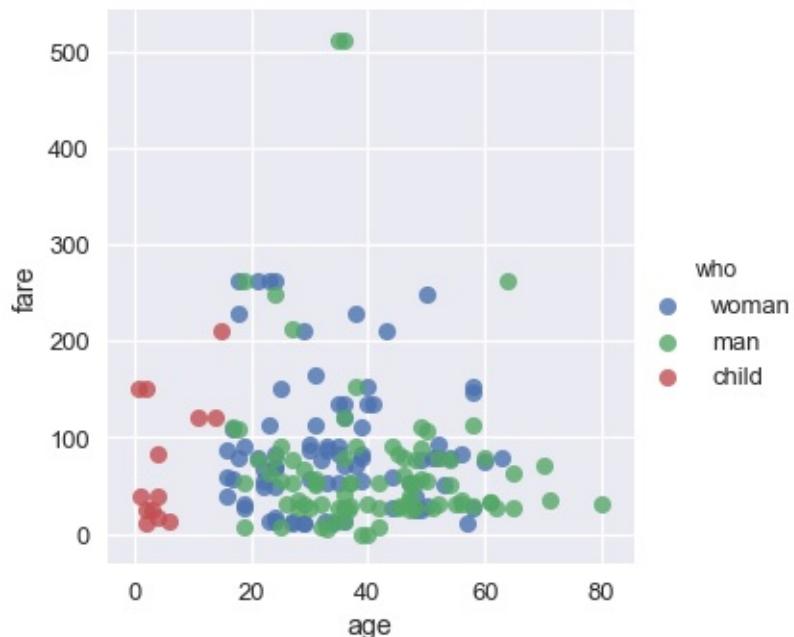
By default seaborn will also fit a regression line to our scatterplot and bootstrap the scatterplot to create a 95% confidence interval around the regression line shown as the light blue shading around the line above. In this case, the regression line doesn't seem to fit the scatter plot very well so we can turn off the regression.

```
sns.lmplot(x='age', y='fare', data=ti, fit_reg=False);
```



We can color the points using a categorical variable. Let's use the `who` column once more:

```
sns.lmplot(x='age', y='fare', hue='who', data=ti,
fit_reg=False);
```



From this plot we can see that all passengers below the age of 18 or so were marked as `child`. There doesn't seem to be a noticeable split between male and female passenger fares, although the two most expensive tickets were purchased by males.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Visualizing Qualitative Data](#)
 - [Bar Charts](#)
 - [Dot Charts](#)

Visualizing Qualitative Data

For qualitative or categorical data, we most often use bar charts and dot charts. We will show how to create these plots using `seaborn` and the Titanic survivors dataset.

```
# Import seaborn and apply its plotting styles
import seaborn as sns
sns.set()

# Load the dataset
ti = sns.load_dataset('titanic').reset_index(drop=True)

# This table is too large to fit onto a page so we'll output
sliders to
# pan through different sections.
df_interact(ti)
```

[Show Widget](#)

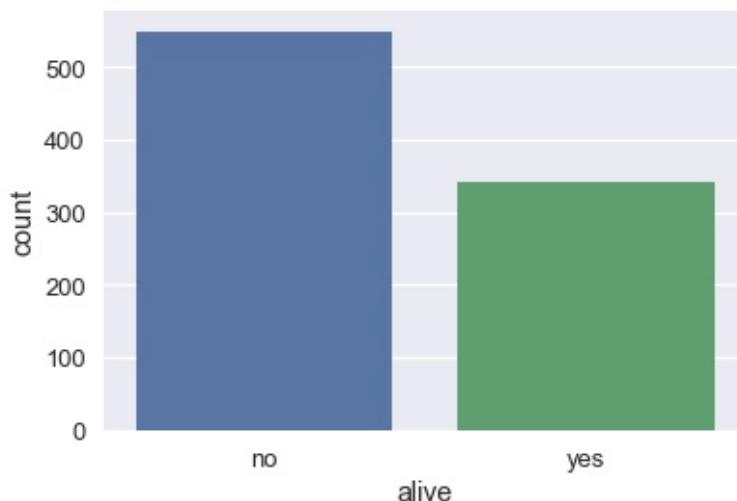
(891 rows, 15 columns) total

Bar Charts

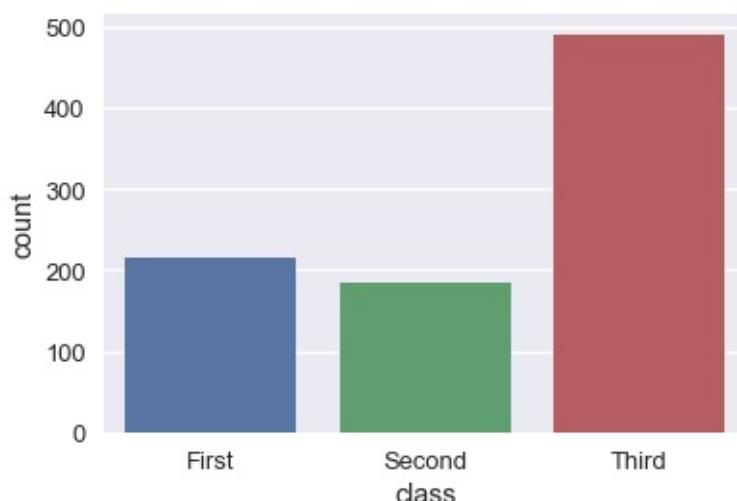
In `seaborn`, there are two types of bar charts. The first type uses the `countplot` method to count up the number of times each category appears in a column.

```
# Counts how many passengers survived and didn't survive and
# draws bars with corresponding heights
sns.countplot(x='alive', data=ti);
```

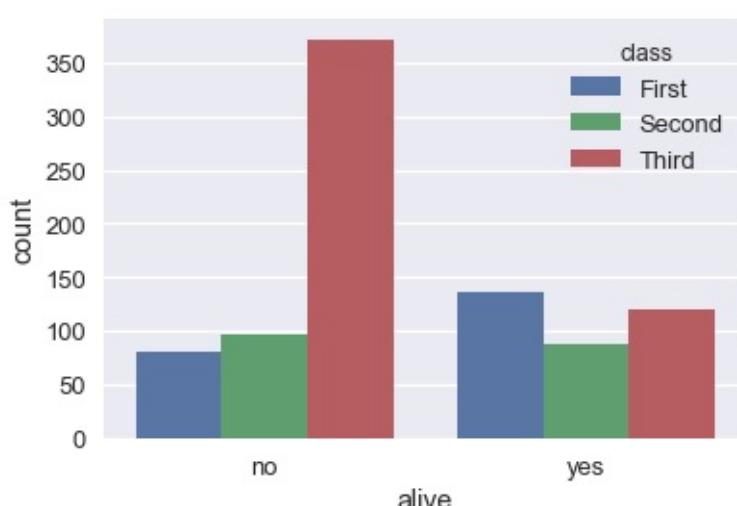
6.2 Qualitative Data



```
sns.countplot(x='class', data=ti);
```

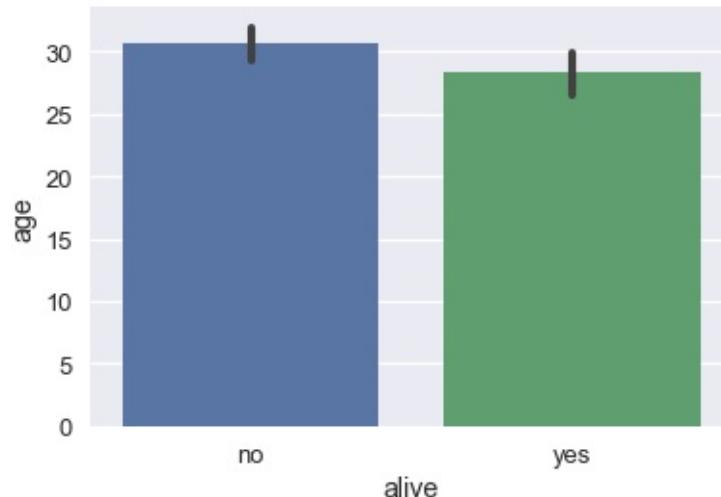


```
# As with box plots, we can break down each category further  
using color  
sns.countplot(x='alive', hue='class', data=ti);
```



The `barplot` method, on the other hand, groups the DataFrame by a categorical column and plots the height of the bars according to the average of a numerical column within each group.

```
# For each set of alive/not alive passengers, compute and plot
# the average age.
sns.barplot(x='alive', y='age', data=ti);
```



The height of each bar can be computed by grouping the original DataFrame and averaging the `age` column:

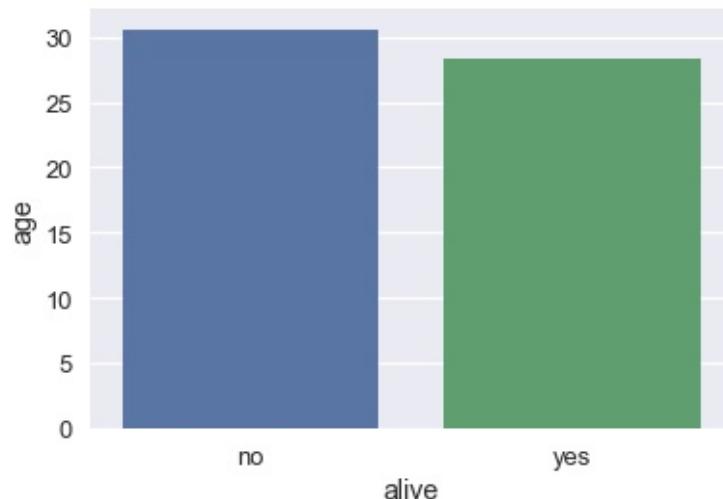
```
ti[['alive', 'age']].groupby('alive').mean()
```

	age
alive	
no	30.626179
yes	28.343690

By default, the `barplot` method will also compute a bootstrap 95% confidence interval for each averaged value, marked as the black lines in the bar chart above. The confidence intervals show that if the dataset contained a random sample of Titanic passengers, the difference between passenger age for those that survived and those that didn't is not statistically significant at the 5% significance level.

These confidence intervals take long to generate when we have larger datasets so it is sometimes useful to turn them off:

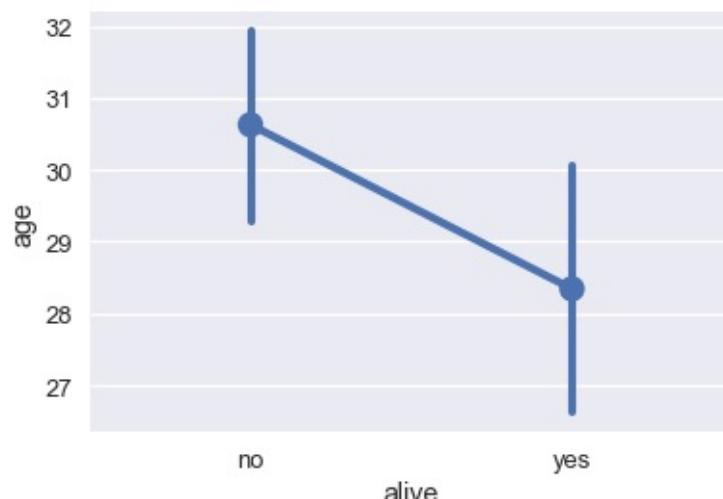
```
sns.barplot(x='alive', y='age', data=ti, ci=False);
```



Dot Charts

Dot charts are similar to bar charts. Instead of plotting bars, dot charts mark a single point at the end of where a bar would go. We use the `pointplot` method to make dot charts in `seaborn`. Like the `barplot` method, the `pointplot` method also automatically groups the DataFrame and computes the average of a separate numerical variable, marking 95% confidence intervals as vertical lines centered on each point.

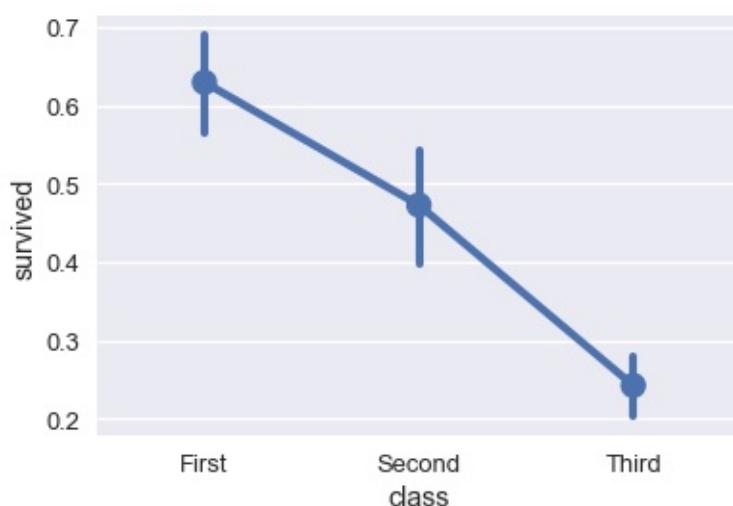
```
# For each set of alive/not alive passengers, compute and plot
# the average age.
sns.pointplot(x='alive', y='age', data=ti);
```



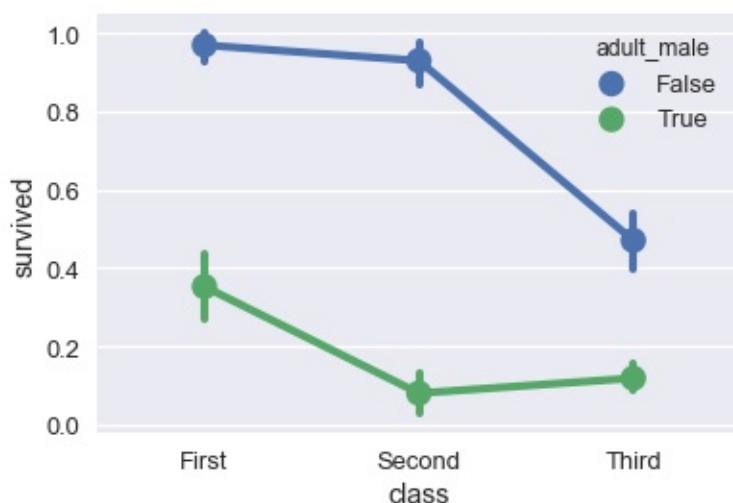
Dot charts are most useful when comparing changes across categories:

```
# Shows the proportion of survivors for each passenger class
sns.pointplot(x='class', y='survived', data=ti);
```

6.2 Qualitative Data



```
# Shows the proportion of survivors for each passenger class,  
# split by whether the passenger was an adult male  
sns.pointplot(x='class', y='survived', hue='adult_male',  
data=ti);
```



[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Customizing Plots using `matplotlib`](#)
 - [Customizing Figures and Axes](#)
 - [Customizing Marks](#)
 - [Arbitrary text and LaTeX support](#)
 - [Customizing a `seaborn` plot using `matplotlib`](#)

Customizing Plots using `matplotlib`

Although `seaborn` allows us to quickly create many types of plots, it does not give us fine-grained control over the chart. For example, we cannot use `seaborn` to modify a plot's title, change x or y-axis labels, or add annotations to a plot. Instead, we must use the `matplotlib` library that `seaborn` is based off of.

`matplotlib` provides basic building blocks for creating plots in Python. Although it gives great control, it is also more verbose—recreating the `seaborn` plots from the previous sections in `matplotlib` would take many lines of code. In fact, we can think of `seaborn` as a set of useful shortcuts to create `matplotlib` plots. Although we prefer to prototype plots in `seaborn`, in order to customize plots for publication we will need to learn basic pieces of `matplotlib`.

Before we look at our first simple example, we must activate `matplotlib` support in the notebook:

```
# This line allows matplotlib plots to appear as images in the
notebook
# instead of in a separate window.
%matplotlib inline

# plt is a commonly used shortcut for matplotlib
import matplotlib.pyplot as plt
```

Customizing Figures and Axes

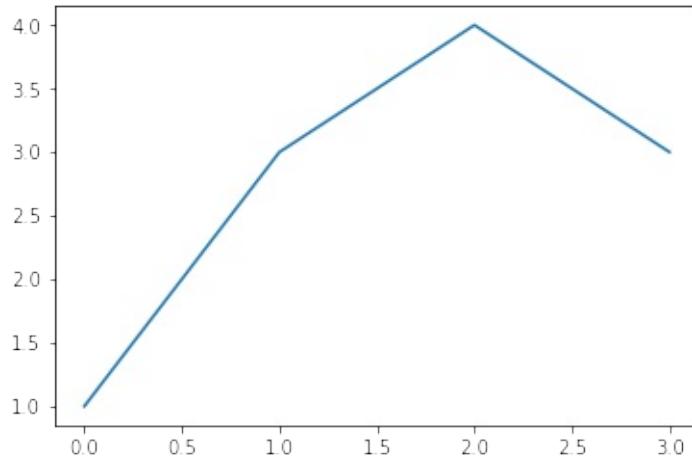
In order to create a plot in `matplotlib`, we create a *figure*, then add an *axes* to the figure. In `matplotlib`, an axes is a single chart, and figures can contain multiple axes in a tabular layout. An axes contains *marks*, the lines or patches drawn on the plot.

```
# Create a figure
f = plt.figure()

# Add an axes to the figure. The second and third arguments
# create a table
# with 1 row and 1 column. The first argument places the axes in
# the first
# cell of the table.
ax = f.add_subplot(1, 1, 1)

# Create a line plot on the axes
ax.plot([0, 1, 2, 3], [1, 3, 4, 3])

# Show the plot. This will automatically get called in a Jupyter
# notebook
# so we'll omit it in future cells
plt.show()
```



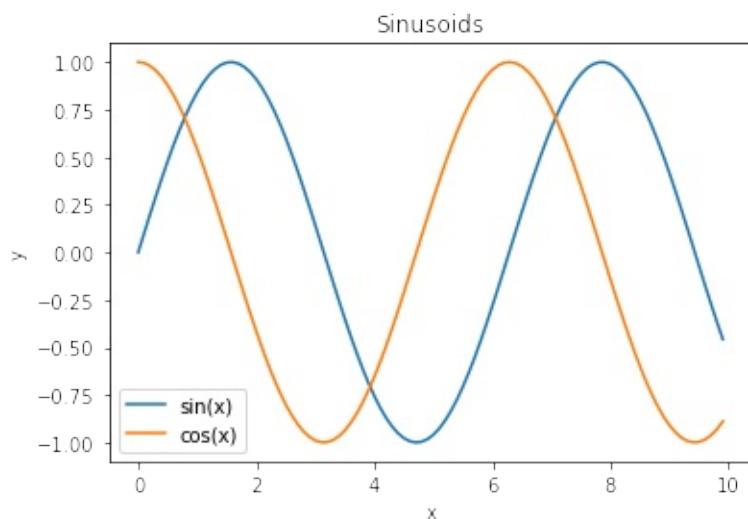
To customize the plot, we can use other methods on the axes object:

```
f = plt.figure()
ax = f.add_subplot(1, 1, 1)

x = np.arange(0, 10, 0.1)

# Setting the label kwarg lets us generate a legend
ax.plot(x, np.sin(x), label='sin(x)')
ax.plot(x, np.cos(x), label='cos(x)')
ax.legend()

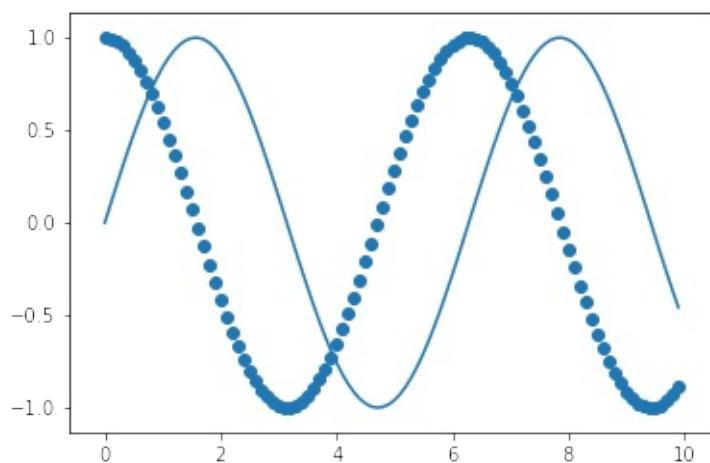
ax.set_title('Sinusoids')
ax.set_xlabel('x')
ax.set_ylabel('y');
```



As a shortcut, `matplotlib` has plotting methods on the `plt` module itself that will automatically initialize a figure and axes.

```
# Shorthand to create figure and axes and call ax.plot
plt.plot(x, np.sin(x))

# When plt methods are called multiple times in the same cell,
# the
# existing figure and axes are reused.
plt.scatter(x, np.cos(x));
```



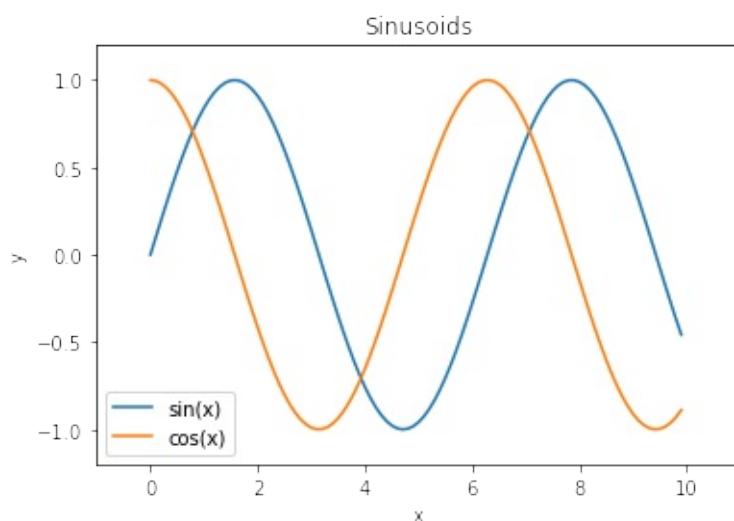
The `plt` module has analogous methods to `an axes`, so we can recreate one of the plots above using `plt` shorthands.

```
x = np.arange(0, 10, 0.1)

plt.plot(x, np.sin(x), label='sin(x)')
plt.plot(x, np.cos(x), label='cos(x)')
plt.legend()

# Shorthand for ax.set_title
plt.title('Sinusoids')
plt.xlabel('x')
plt.ylabel('y')

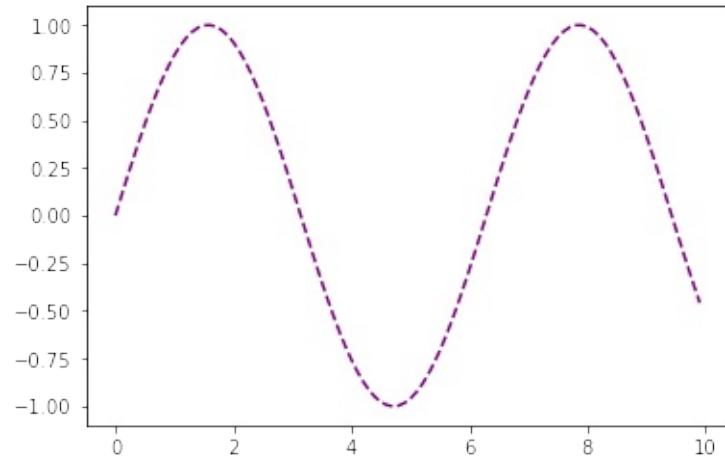
# Set the x and y-axis limits
plt.xlim(-1, 11)
plt.ylim(-1.2, 1.2);
```



Customizing Marks

To change properties of the plot marks themselves (e.g. the lines in the plot above), we can pass additional arguments into `plt.plot`.

```
plt.plot(x, np.sin(x), linestyle='--', color='purple');
```



Checking the `matplotlib` documentation is the easiest way to figure out which arguments are available for each method. Another way is to store the returned line object:

```
In [1]: line, = plot([1,2,3])
```

These line objects have a lot of properties you can control, here's the full list using tab-completion in IPython:

```
In [2]: line.set
line.set           line.set_drawstyle      line.set_mec
line.set_aa        line.set_figure        line.set_mew
line.set_agg_filter line.set_fillstyle    line.set_mfc
line.set_alpha     line.set_gid          line.set_mfcalt
line.set_animated   line.set_label        line.set_ms
line.set_antialiased line.set_linestyle   line.set_picker
line.set_axes       line.set_linewidth    line.set_pickradius
line.set_c          line.set_lod          line.set_rasterized
line.set_clip_box   line.set_ls           line.set_snap
line.set_clip_on    line.set_lw           line.set_solid_capstyle
line.set_clip_path  line.set_marker       line.set_solid_joinstyle
line.set_color      line.set_markeredgewidth line.set_transform
line.set_contains   line.set_markeredgecolor line.set_url
line.set_dash_capstyle line.set_markerfacecolor line.set_visible
line.set_dashes     line.set_markerfacecoloralt line.set_xdata
line.set_dash_joinstyle line.set_markersize   line.set_ydata
line.set_data       line.set_markevery    line.set_zorder
```

But the `setp` call (short for `set property`) can be very useful, especially while working interactively because it contains introspection support, so you can learn about the valid calls as you work:

```
In [7]: line, = plot([1,2,3])

In [8]: setp(line, 'linestyle')
linestyle: [ '-' | '--' | '-.' | ':' | 'None' | ' ' | ' ' ]
]       and any drawstyle in combination with a linestyle, e.g. ``'steps--'``.

In [9]: setp(line)
agg_filter: unknown
alpha: float (0.0 transparent through 1.0 opaque)
animated: [True | False]
antialiased or aa: [True | False]
...
... much more output omitted
...
```

In the first form, it shows you the valid values for the 'linestyle' property, and in the second it shows you all the acceptable properties you can set on the line object. This makes it easy to discover how to customize your figures to get the visual results you need.

Arbitrary text and LaTeX support

In matplotlib, text can be added either relative to an individual axis object or to the whole figure.

These commands add text to the Axes:

- `set_title()` - add a title
- `set_xlabel()` - add an axis label to the x-axis
- `set_ylabel()` - add an axis label to the y-axis
- `text()` - add text at an arbitrary location
- `annotate()` - add an annotation, with optional arrow

And these act on the whole figure:

- `figtext()` - add text at an arbitrary location
- `suptitle()` - add a title

And any text field can contain LaTeX expressions for mathematics, as long as they are enclosed in `$` signs.

This example illustrates all of them:

```
fig = plt.figure()
fig.suptitle('bold figure suptitle', fontsize=14,
fontweight='bold')

ax = fig.add_subplot(1, 1, 1)
fig.subplots_adjust(top=0.85)
ax.set_title('axes title')

ax.set_xlabel('xlabel')
ax.set_ylabel('ylabel')

ax.text(3, 8, 'boxed italics text in data coords',
style='italic',
bbox={'facecolor':'red', 'alpha':0.5, 'pad':10})

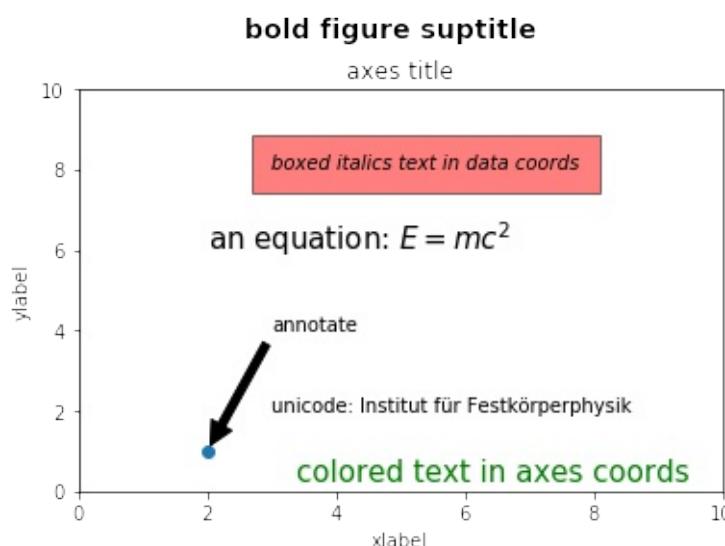
ax.text(2, 6, 'an equation: $E=mc^2$', fontsize=15)

ax.text(3, 2, 'unicode: Institut für Festkörperphysik')

ax.text(0.95, 0.01, 'colored text in axes coords',
verticalalignment='bottom', horizontalalignment='right',
transform=ax.transAxes,
color='green', fontsize=15)

ax.plot([2], [1], 'o')
ax.annotate('annotate', xy=(2, 1), xytext=(3, 4),
arrowprops=dict(facecolor='black', shrink=0.05))

ax.axis([0, 10, 0, 10]);
```



Customizing a `seaborn` plot using `matplotlib`

Now that we've seen how to use `matplotlib` to customize a plot, we can use the same methods to customize `seaborn` plots since `seaborn` creates plots using `matplotlib` behind-the-scenes.

```
# Load seaborn
import seaborn as sns
sns.set()
sns.set_context('talk')

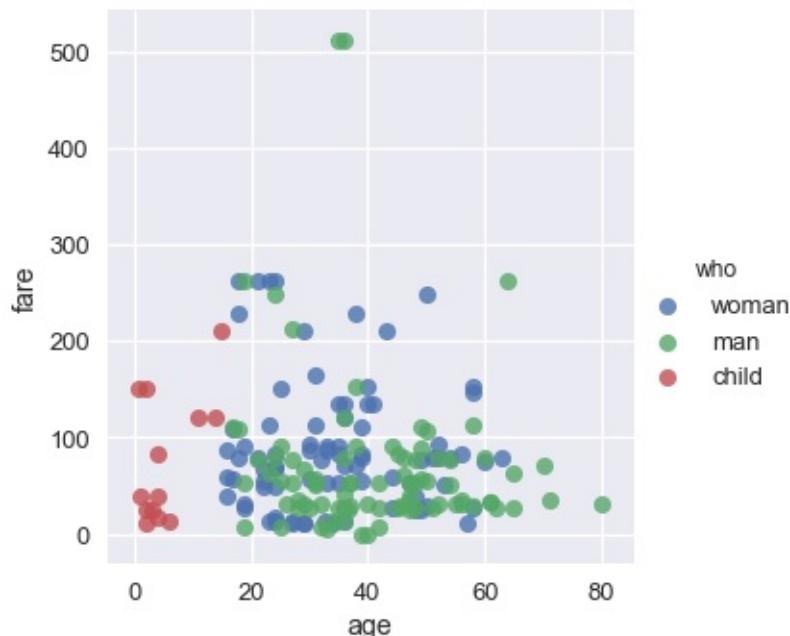
# Load dataset
ti = sns.load_dataset('titanic').dropna().reset_index(drop=True)
ti.head()
```

	survived	pclass	sex	age	...	deck	embark_town	alive	...
0	1	1	female	38.0	...	C	Cherbourg	yes	1
1	1	1	female	35.0	...	C	Southampton	yes	1
2	0	1	male	54.0	...	E	Southampton	no	1
3	1	3	female	4.0	...	G	Southampton	yes	1
4	1	1	female	58.0	...	C	Southampton	yes	1

5 rows × 15 columns

We'll start with this plot:

```
sns.lmplot(x='age', y='fare', hue='who', data=ti,
fit_reg=False);
```



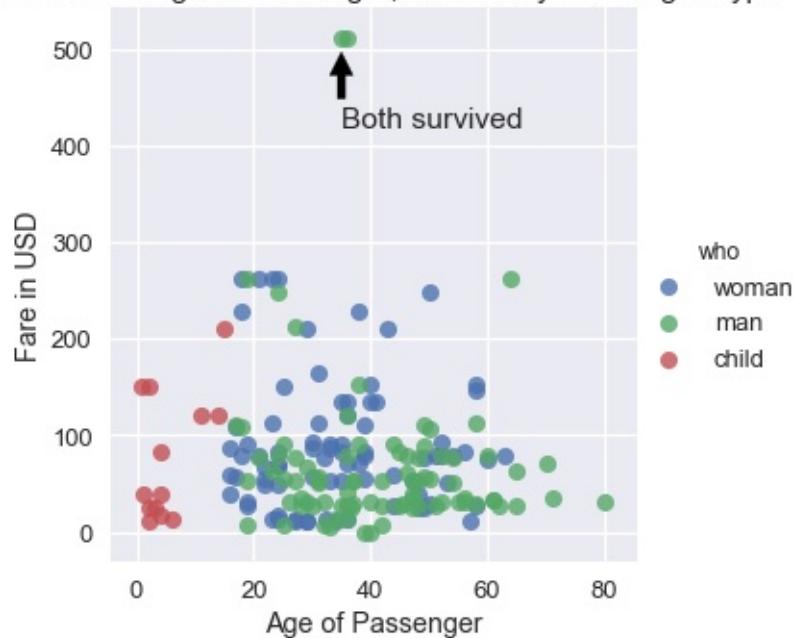
We can see that the plot needs a title and better labels for the x and y-axes. In addition, the two people with the most expensive fares survived, so we can annotate them on our plot.

```
sns.lmplot(x='age', y='fare', hue='who', data=ti, fit_reg=False)

plt.title('Fare Paid vs. Age of Passenger, Colored by Passenger
Type')
plt.xlabel('Age of Passenger')
plt.ylabel('Fare in USD')

plt.annotate('Both survived', xy=(35, 500), xytext=(35, 420),
arrowprops=dict(facecolor='black', shrink=0.05));
```

Fare Paid vs. Age of Passenger, Colored by Passenger Type



In practice, we use `seaborn` to quickly explore the data and then turn to `matplotlib` for fine-tuning once we decide on the plots to use in a paper or presentation.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Visualization Principles](#)
 - [Principles of Scale](#)
 - [Principles of Conditioning](#)
 - [Principles of Perception](#)

Visualization Principles¶

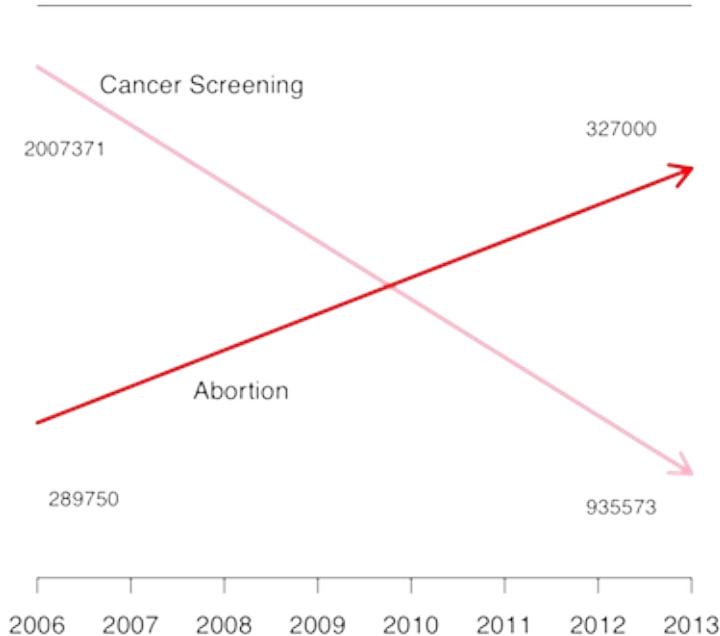
Now that we have the tools to create and alter plots, we turn to key principles for data visualization. Much like other parts of data science, it is difficult to precisely assign a number that measures how effective a specific visualization is. Still, there are general principles that make visualizations much more effective at showing trends in the data. We discuss six categories of principles: scale, conditioning, perception, transformation, context, and smoothing.

Principles of Scale¶

Principles of scale relate to the choice of x and y-axis used to plot the data.

In a 2015 US Congressional hearing, representative Chaffetz discussed an investigation of Planned Parenthood programs. He presented the following plot that originally appeared in a report by Americans United for Life. It compares the number of abortion and cancer screening procedures, both of which are offered by Planned Parenthood. (The full report is available at <https://oversight.house.gov/interactivepage/plannedparenthood>.)

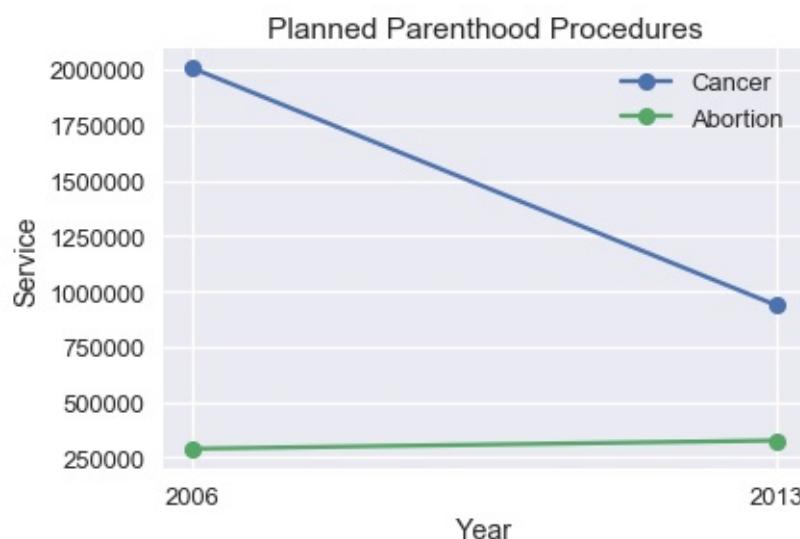
What is suspicious about this plot? How many data points are plotted?



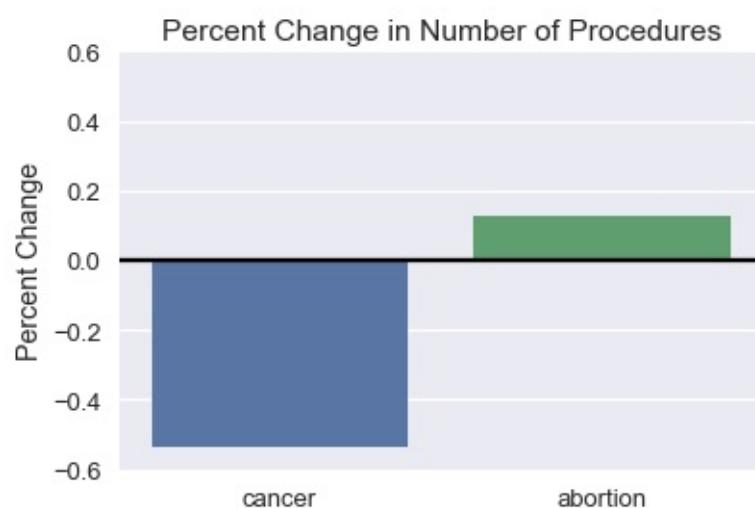
This plot violates principles of scale; it doesn't make good choices for its x and y-axis.

When we select the x and y-axis for our plot, we should keep a consistent scale across the entire axis. However, the plot above has different scales for the Abortion and Cancer Screening lines—the start of the Abortion line and end of the Cancer Screening line lie close to each other on the y-axis but represent vastly different numbers. In addition, only points from 2006 and 2013 are plotted but the x-axis contains unnecessary tick marks for every year in between.

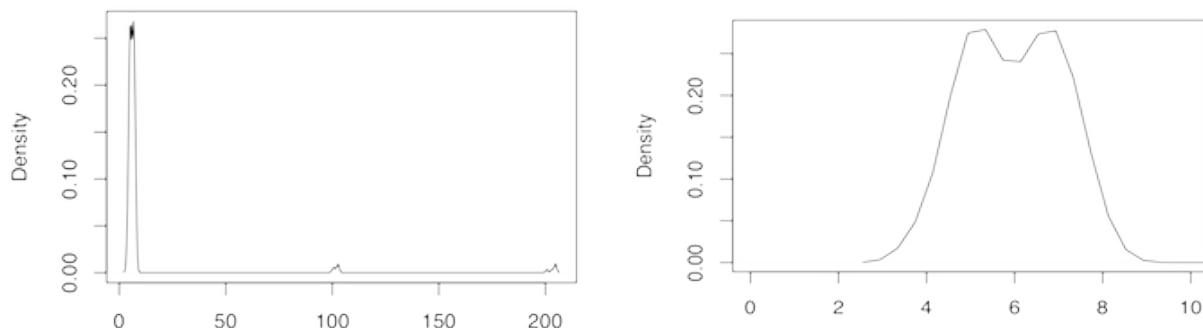
To improve this plot, we should re-plot the points on the same y-axis scale:



We can see that the change in number of Abortions is very small compared to the large drop in the number of Cancer Screenings. Instead of the number of procedures, we might instead be interested in the percent change in number.



When selecting the x and y-axis limits we prefer to focus on the region with the bulk of the data, especially when working with long-tailed data. Consider the following plot and its zoomed in version to its right:



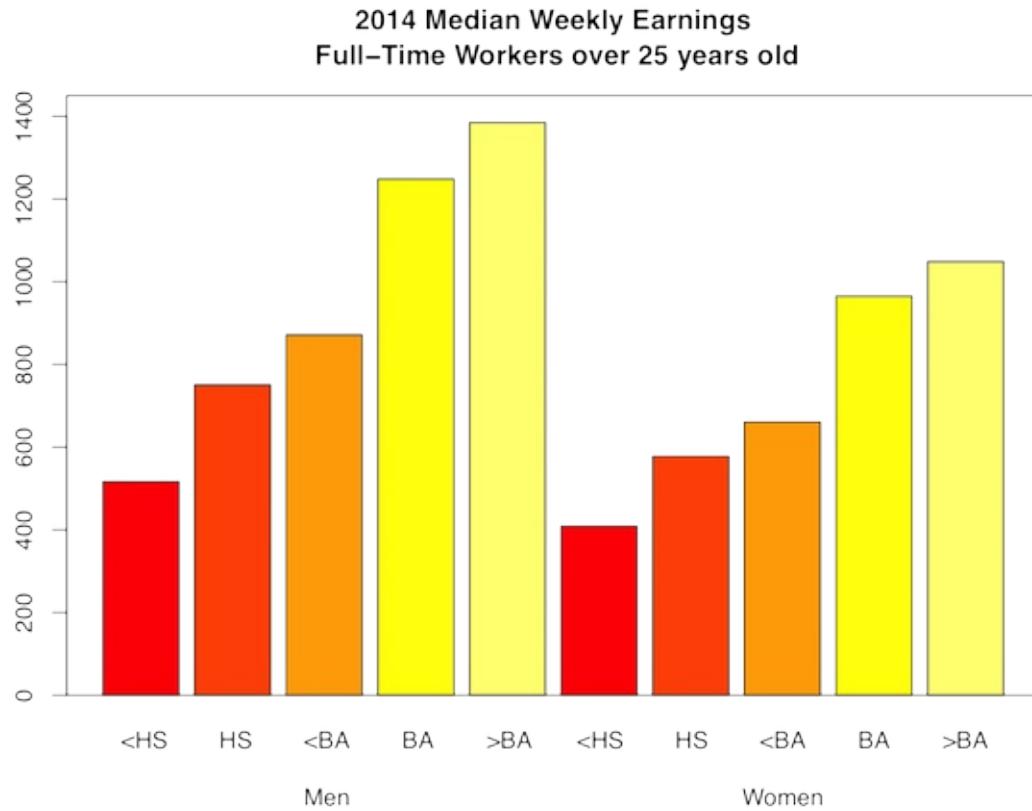
The plot on the right is much more helpful for making sense of the dataset. If needed, we can make multiple plots of different regions of the data to show the entire range of data. Later in this section, we discuss data transformations which also help visualize long-tailed data.

Principles of Conditioning

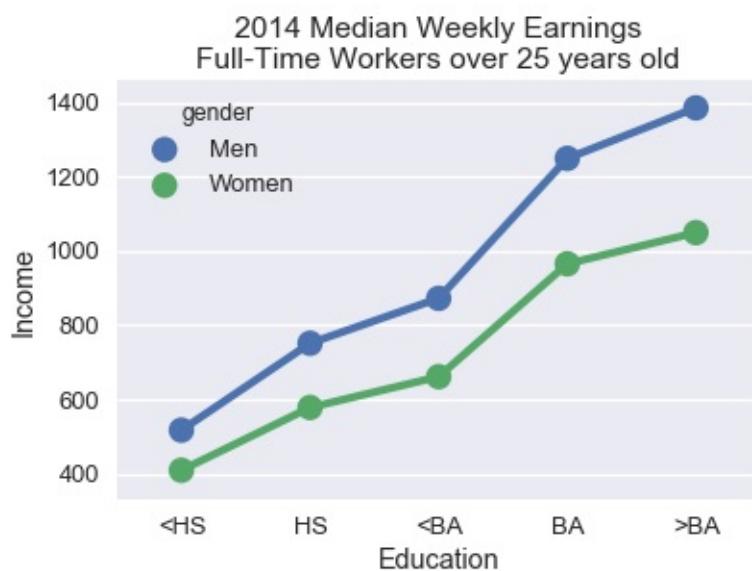
Principles of conditioning give us techniques to show distributions and relationships between subgroups of our data.

The US Bureau of Labor Statistics oversees scientific surveys related to the economic health of the US. Their website contains a tool to generate reports using this data that was used to generate this chart comparing median weekly earnings split by sex.

Which comparisons are easiest to make using this plot? Are these the comparisons that are most interesting or important?



This plot lets us see at a glance that weekly earnings tend to increase with more education. However, it is difficult to tell exactly how much each level of education increases earnings and it is even more difficult to compare male and female weekly earnings at the same education level. We can uncover both these trends by using a dot chart instead of a bar chart.

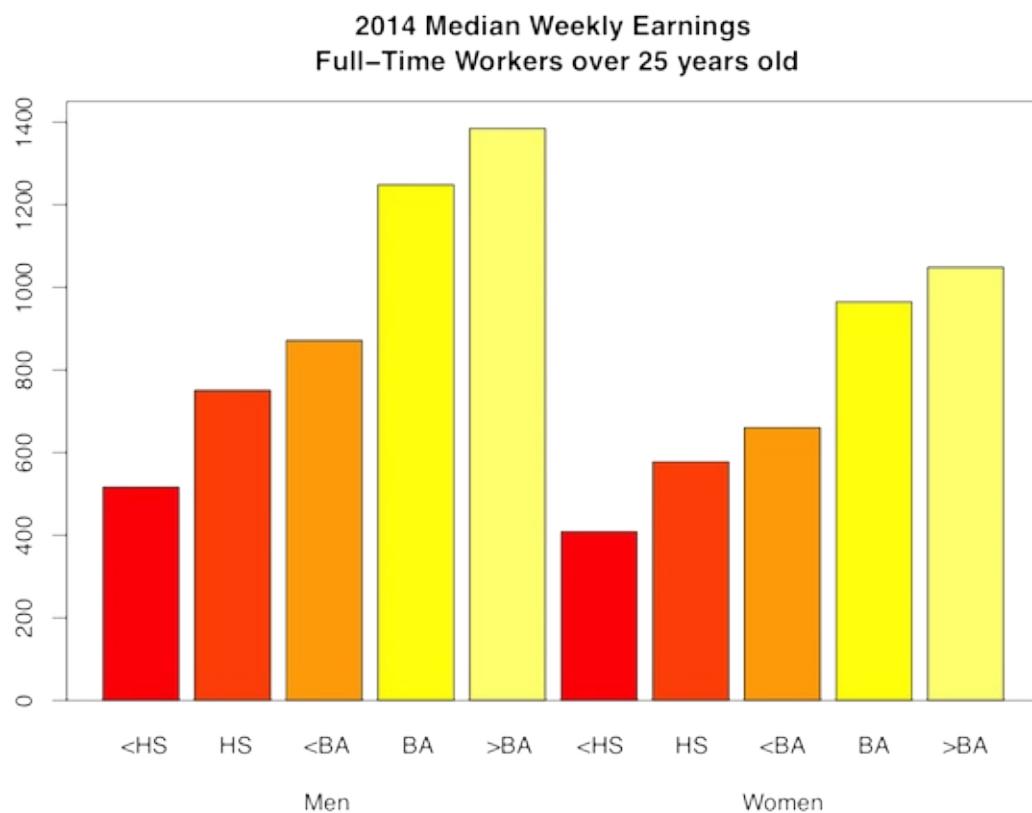


The lines connecting the points more clearly shows the relatively large effect of having a BA degree on weekly earnings. Placing the points for males and females directly above each other makes it much easier to see that the wage gap between males and females tends to increase with higher education levels.

To aid comparison of two subgroups within your data, align markers along the x or y-axis and use different colors or markers for different subgroups. Lines tend to show trends in data more clearly than bars and are a useful choice for both ordinal and numerical data.

Principles of Perception

Human perception has specific properties that are important to consider in visualization design. The first important property of human perception is that we perceive some colors more strongly than others, especially green colors. In addition, we perceive lighter shaded areas as larger than darker shaded ones. For example, in the weekly earnings plot that we just discussed, the lighter bars seem to draw more attention than the darker colored ones:



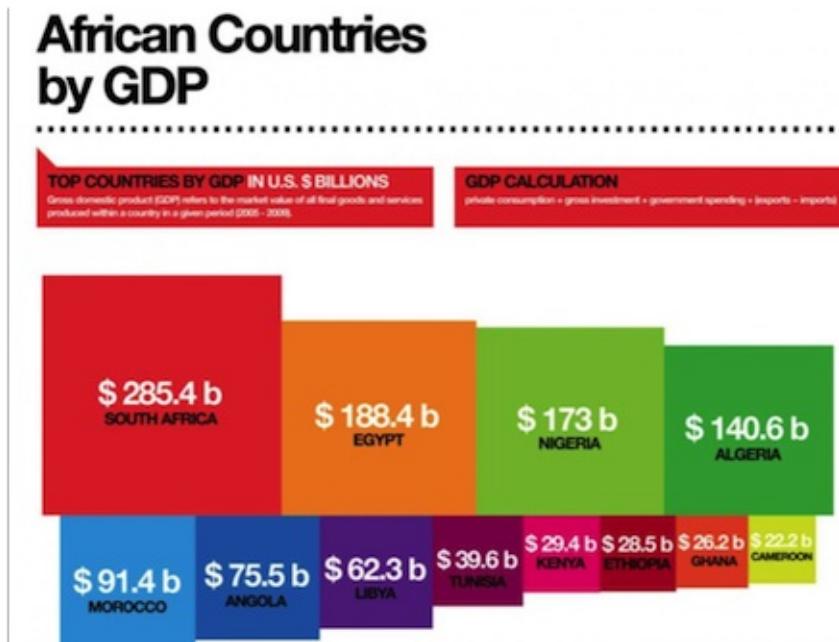
Practically speaking, you should ensure that your charts' color palettes are *perceptually uniform*. This means that, for example, the perceived intensity of the color won't change in between bars in a bar chart. For quantitative data, you have two choices: if your data progress from low to high and you want to emphasize large values, use a *sequential* color

scheme which assigns lighter colors to large values. If both low and high values should be emphasized, use a *diverging* color scheme which assigns lighter colors to values closer to the center.

`seaborn` comes with many useful color palettes built-in. You can browse its documentation to learn how to switch between color palettes:

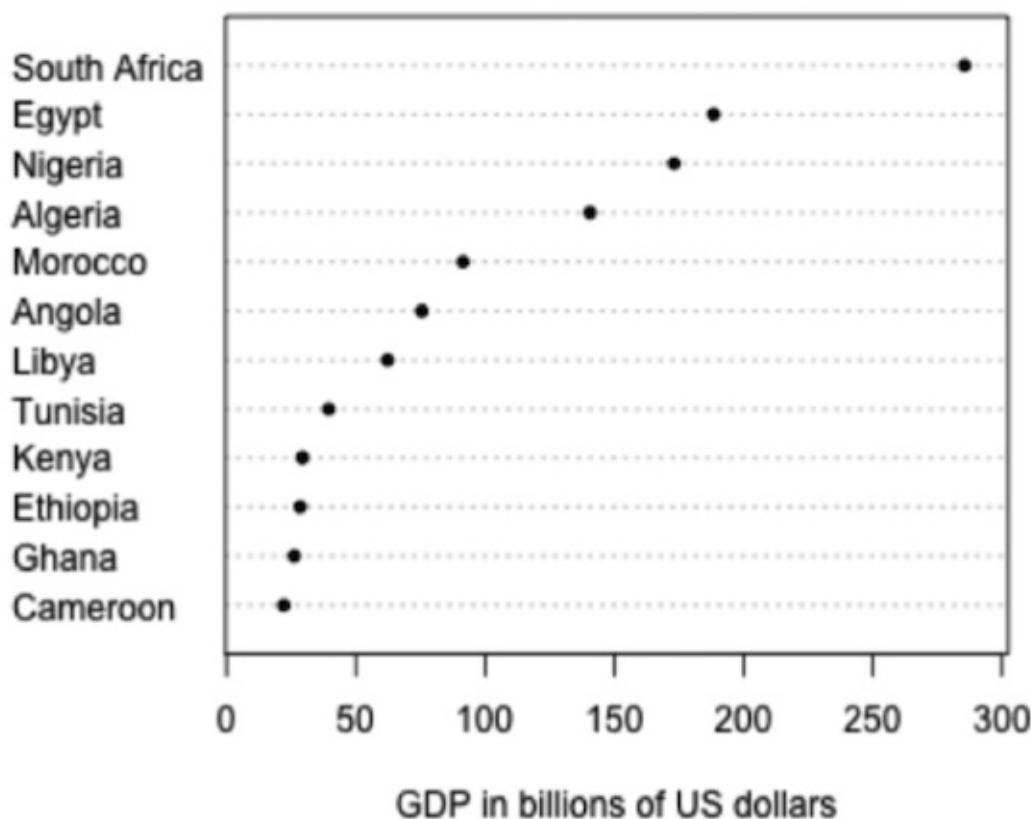
http://seaborn.pydata.org/tutorial/color_palettes.html

A second important property of human perception is that we are generally more accurate when we compare lengths and less accurate when we compare areas. Consider the following chart of the GDP of African countries.



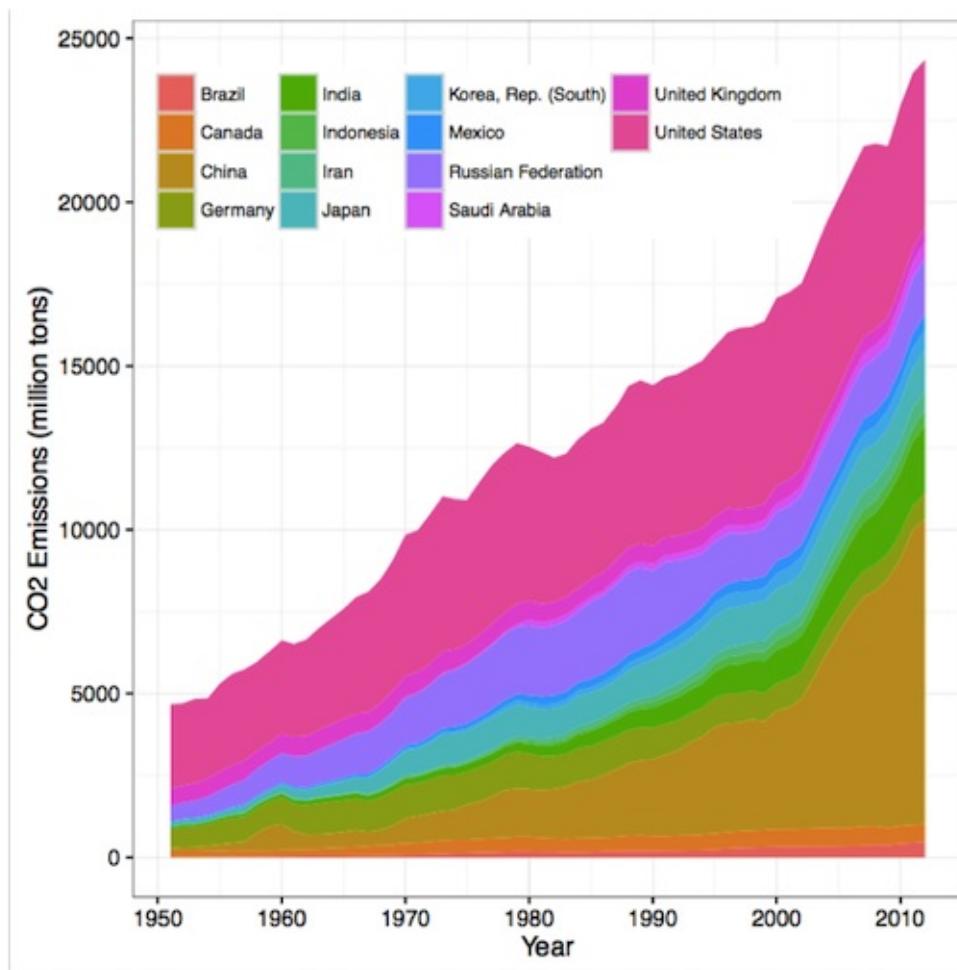
By numerical value, South Africa has twice the GDP of Algeria but it's not easy to tell from the plot above. Instead, we can plot the GDPs on a dot plot:

African Countries by GDP



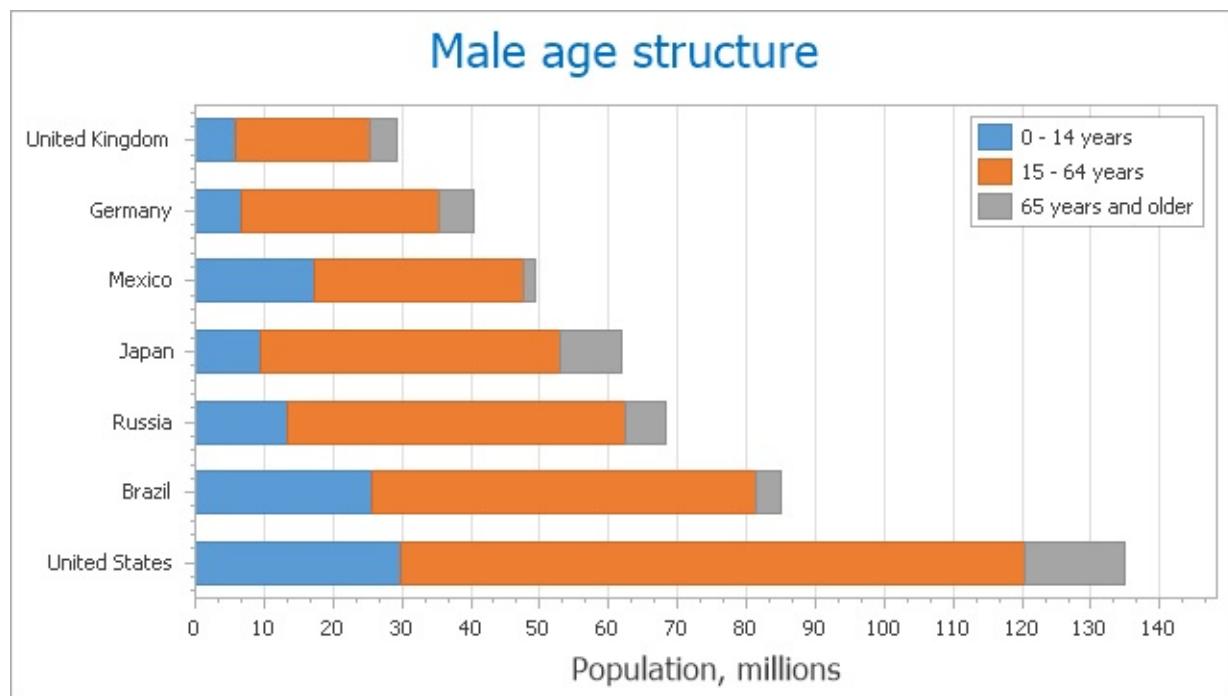
This is much more clear because it allows us to compare lengths instead of areas. Pie charts and three-dimensional charts are difficult to interpret for the same reason; we tend to avoid these charts in practice.

Our third and final property of perception is that the human eye has difficulty with changing baselines. Consider the following stacked area chart that plots carbon dioxide emissions over time split by country.

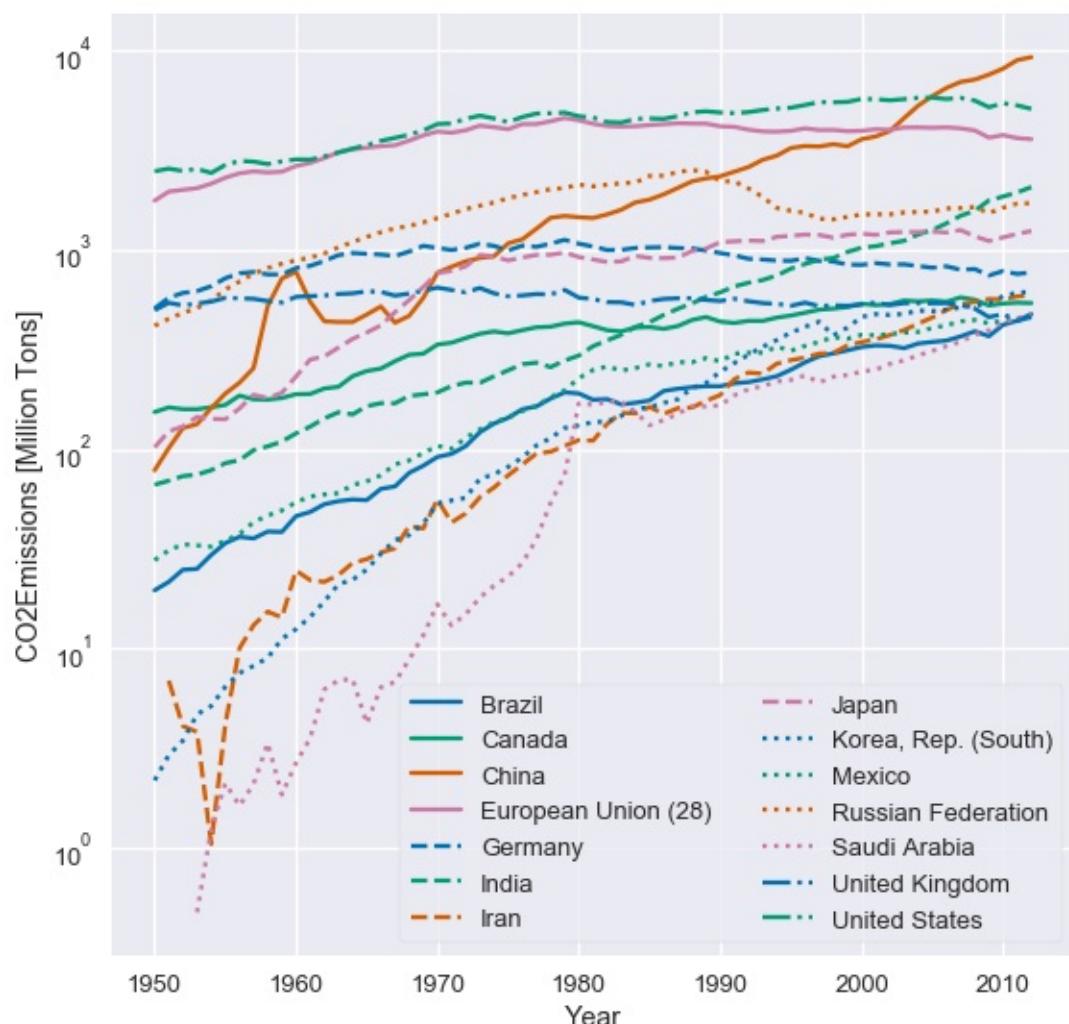


It is difficult to see whether the UK's emissions have increased or decreased over time because of the *jiggling baseline* problem: the baseline (bottom line) of the area jiggles up and down. It is also difficult to compare whether the UK's emissions are greater than China's emissions when the two heights are similar (in year 2000, for example).

Similar issues of jiggling baselines appear in stacked bar charts. In the plot below, it is difficult to compare the number of 15-64 year olds between Germany and Mexico.



We can often improve a stacked area or bar chart by switching to a line chart. Here's the data of emissions over time plotted as lines instead of areas:



This plot does not jiggle the baseline so it is much easier to compare emissions between countries. We can also more clearly see which countries increased emissions the most.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Visualization Principles Continued](#)
 - [Principles of Transformation](#)
 - [Principles of Context](#)
 - [Principles of Smoothing](#)
 - [Kernel Density Estimation Details](#)
 - [Smoothing a Scatter Plot](#)

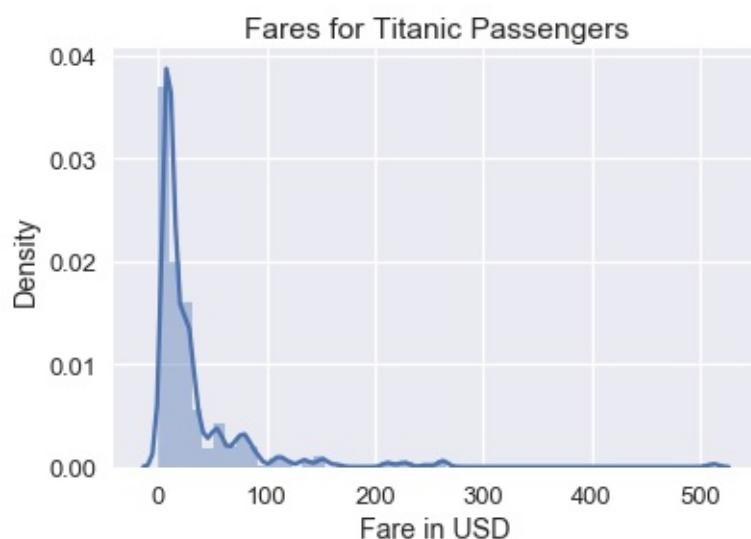
Visualization Principles Continued

In this section, we discuss principles of visualization for transformation, context, and smoothing.

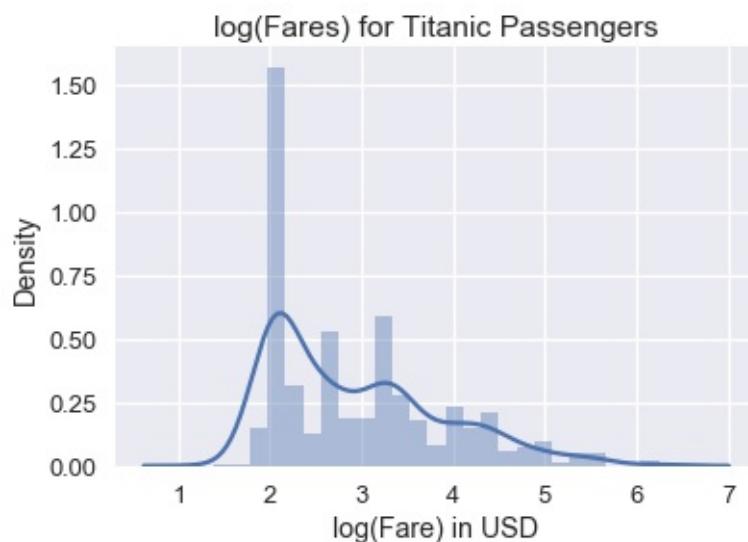
Principles of Transformation

The principles of data transformation give us useful ways to alter data for visualization in order to more effectively reveal trends. We most commonly apply data transformations to reveal patterns in skewed data and non-linear relationships between variables.

The plot below shows the distribution of ticket fares for each passenger aboard the Titanic. As you can see, the distribution is skewed right.



Although this histogram shows all the fares, it is difficult to see detailed patterns in the data since the fares are clumped on the left side of the histogram. To remedy this, we can take the natural log of the fares before plotting them:



We can see from the plot of the log data that the distribution of fares has a mode at roughly $\$ e^2 = \7.40 and a smaller mode at roughly $\$ e^{3.4} = \30.00 . Why does plotting the natural log of the data help with skew? The logarithms of large numbers tend to be close to the logarithms of small numbers:

value	log(value)
1	0.00
10	2.30
50	3.91
100	4.60
500	6.21
1000	6.90

This means that taking the logarithm of right-tailed data will bring large values close to small values. This helps see patterns where the majority of the data lie.

In fact, the logarithm is considered the Swiss army knife of data transformation—it also helps us see the nature of non-linear relationships between variables in the data. In 1619, Kepler recorded down the following set of data to discover his Third Law of Planetary Motion:

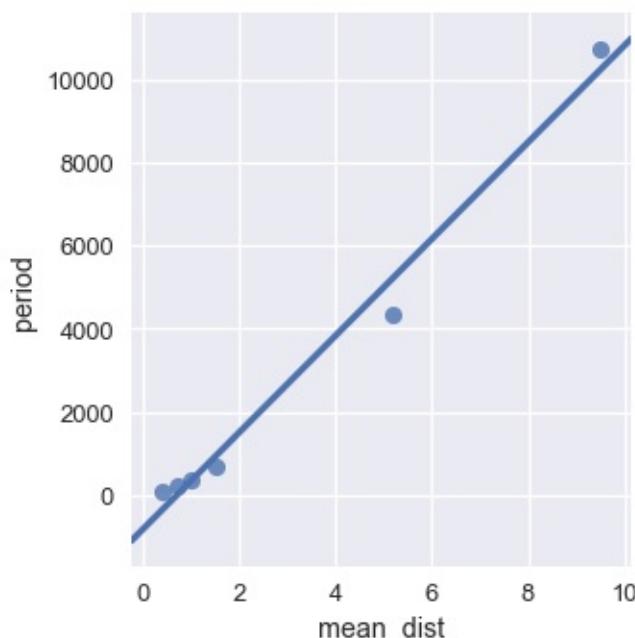
```
planets = pd.read_csv("data/planets.data",
                      delim_whitespace=True,
                      comment="#", usecols=[0, 1, 2])
planets
```

	planet	mean_dist	period
0	Mercury	0.389	87.77
1	Venus	0.724	224.70
2	Earth	1.000	365.25
3	Mars	1.524	686.95
4	Jupiter	5.200	4332.62
5	Saturn	9.510	10759.20

If we plot the mean distance to the sun against the period of the orbit, we can see a relationship that doesn't quite look linear:

```
sns.lmplot(x='mean_dist', y='period', data=planets, ci=False)
```

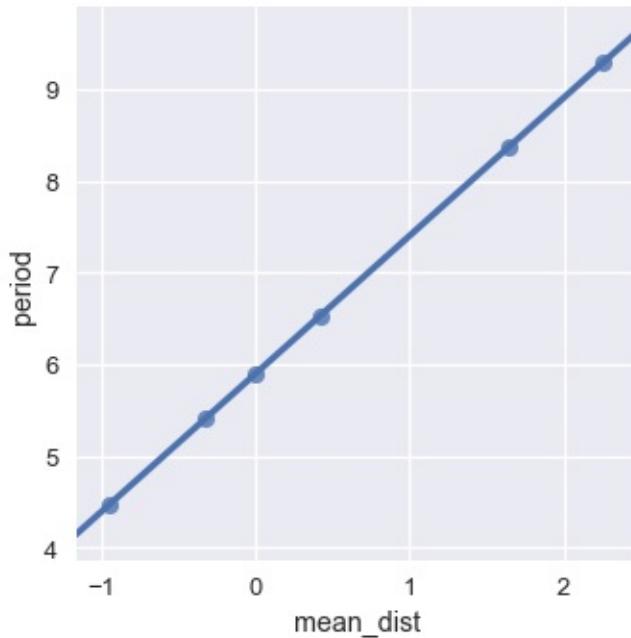
```
<seaborn.axisgrid.FacetGrid at 0x1a1f54aba8>
```



However, if we take the natural log of both mean distance and period, we obtain the following plot:

```
sns.lmplot(x='mean_dist', y='period',
            data=np.log(planets.iloc[:, [1, 2]]),
            ci=False);
```

```
<seaborn.axisgrid.FacetGrid at 0x1a1f693da0>
```



We see a near-perfect linear relationship between the logged values of mean distance and period. What does this mean? Since we believe there's a linear relationship between the logged values, we can derive:

$$\begin{aligned}
 \log(\text{period}) &= m \log(\text{dist}) + b \\
 \text{Taking the exponent of both sides} \quad \text{period} &= e^b \text{dist}^m \\
 \text{period} &= C \cdot \text{dist}^m
 \end{aligned}$$

We replaced e^b with C in the last step to represent e^b as a constant. The algebraic manipulation above shows that when two variables have a polynomial relationship, the log of the two variables has a linear relationship. In fact, we can find the degree of the polynomial by examining the slope of the line. In this case, the slope is 1.5 which gives us Kepler's third law: $\text{period} \propto \text{dist}^{1.5}$.

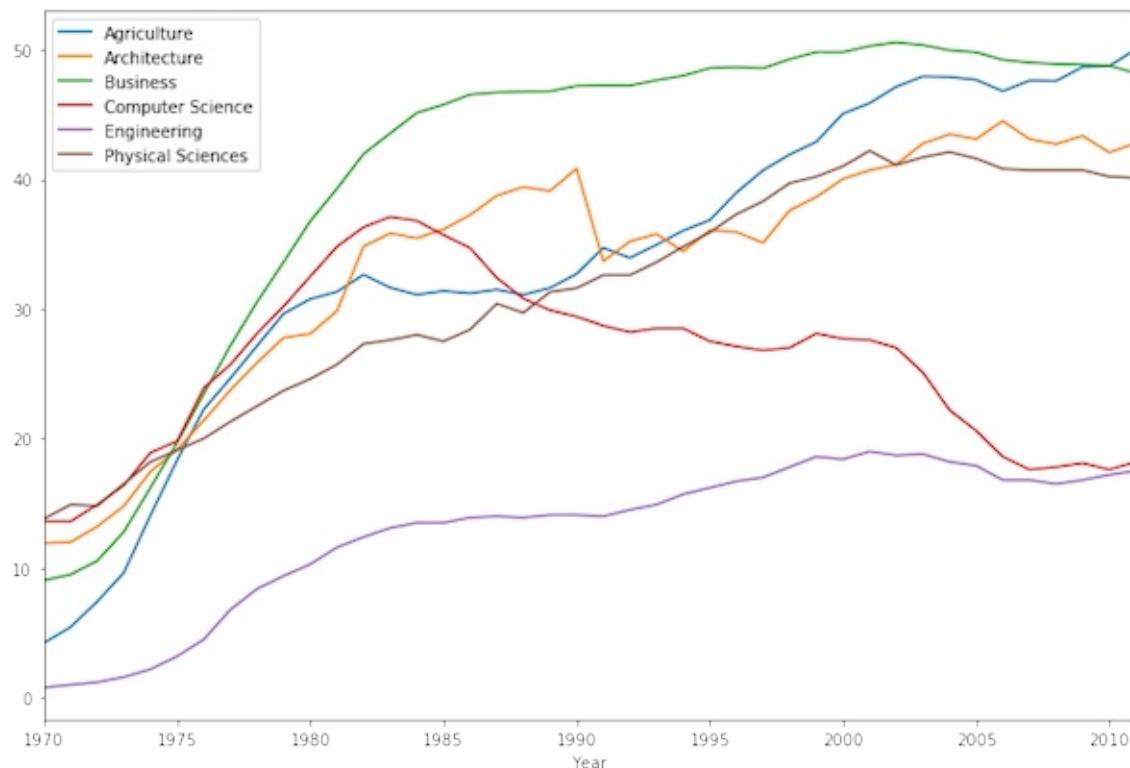
By a similar derivation we can also show that if the relationship between $\log(y)$ and x is linear, the two variables have an exponential relationship: $y = a^x$.

Thus, we can use the logarithm to reveal patterns in right-tailed data and common non-linear relationships between variables.

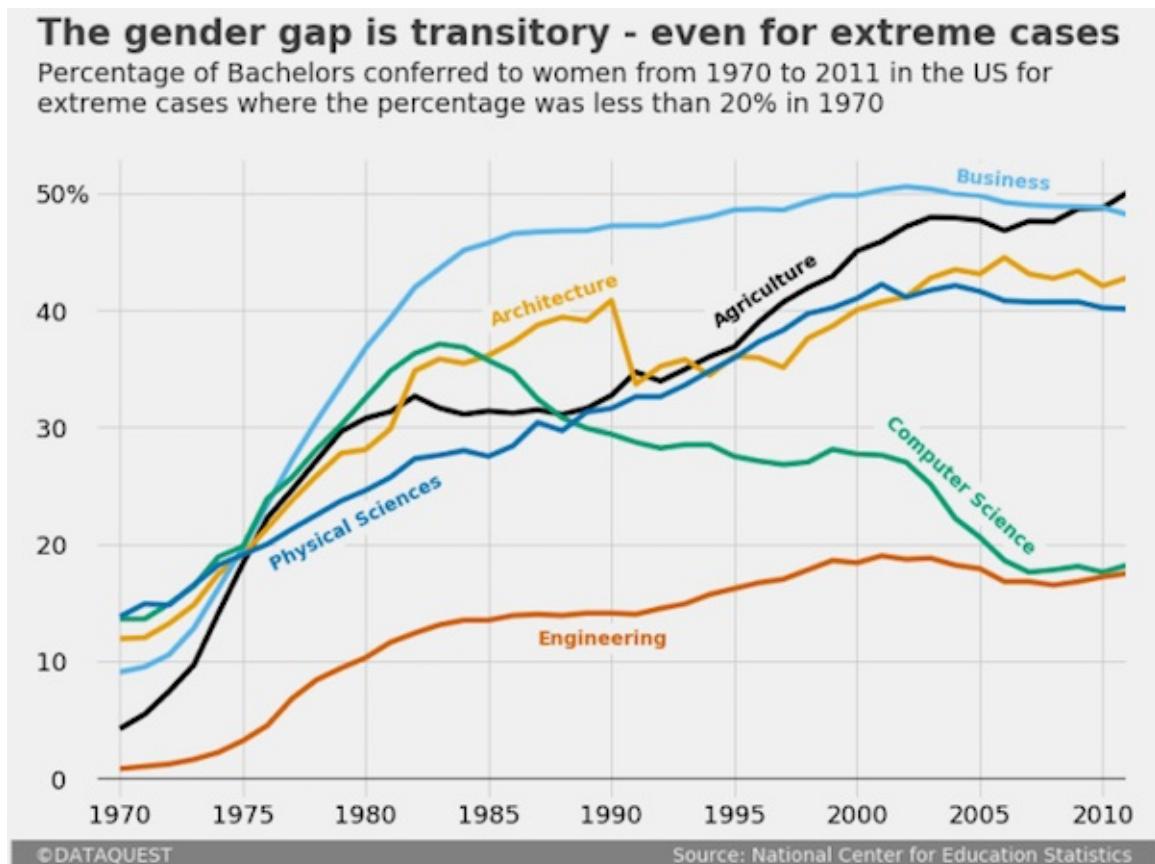
Other common data transformations include the Box-Cox transformation and polynomial transforms.

Principles of Context

It is important to add as much relevant context as possible to any plot you plan to share more broadly. For example, the following plot shows its data clearly but provides little context to help understand what is being plotted.



To provide context, we add a title, caption, axes labels, units for the axes, and labels for the plotted lines.



([This blog post](#) explains how to make these modifications using `matplotlib`.)

In general, we provide context for a plot through:

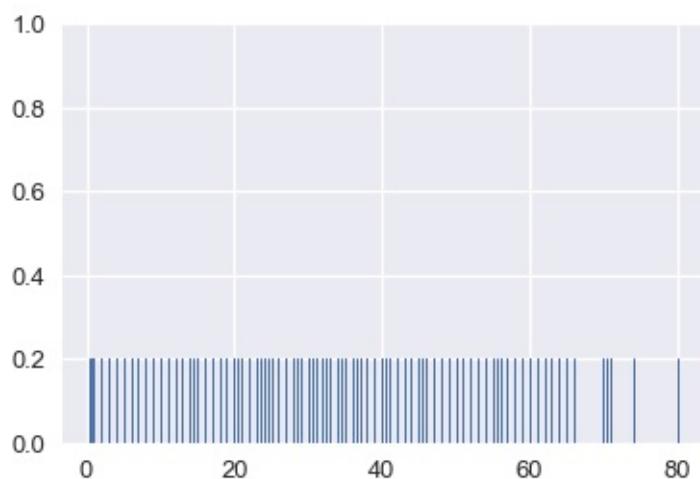
- Plot title
- Axes labels
- Reference lines and markers for important values
- Labels for interesting points and unusual observations
- Captions that describe the data and its important features

Principles of Smoothing

Smoothing allows us to more clearly visualize data when we have many data points. We've actually already seen an instance of smoothing: histograms are a type of smoothing for rugplots. This rugplot shows each age of the passengers in the Titanic.

```
ages = ti['age'].dropna()
sns.rugplot(ages, height=0.2)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a20c05b38>
```

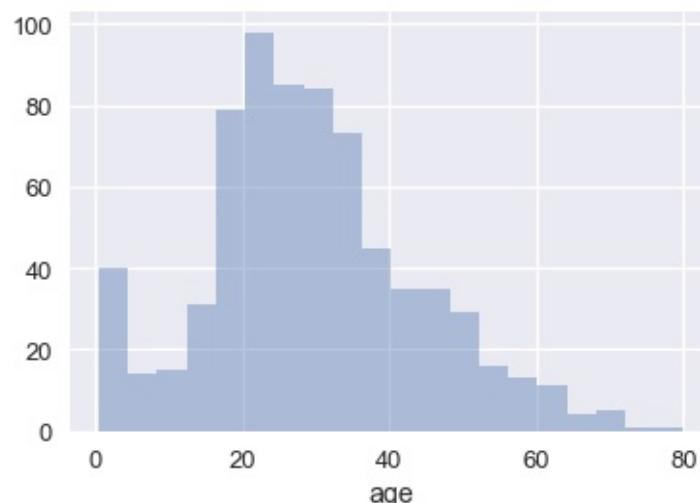


There are many marks that make it difficult to tell where the data lie. In addition, some of the points overlap, making it impossible to see how many points lie at 0. This issue is called *overplotting* and we generally avoid it whenever possible.

To reveal the distribution of the data, we can replace groups of marks with a bar that is taller when more points are in the group. Smoothing refers to this process of replacing sets of points with appropriate markers; we choose not to show every single point in the dataset in order to reveal broader trends.

```
sns.distplot(ages, kde=False)
```

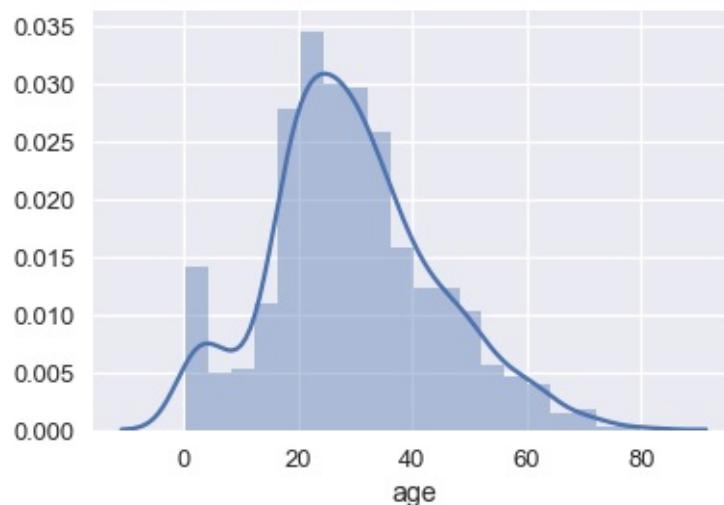
```
<matplotlib.axes._subplots.AxesSubplot at 0x1a23c384e0>
```



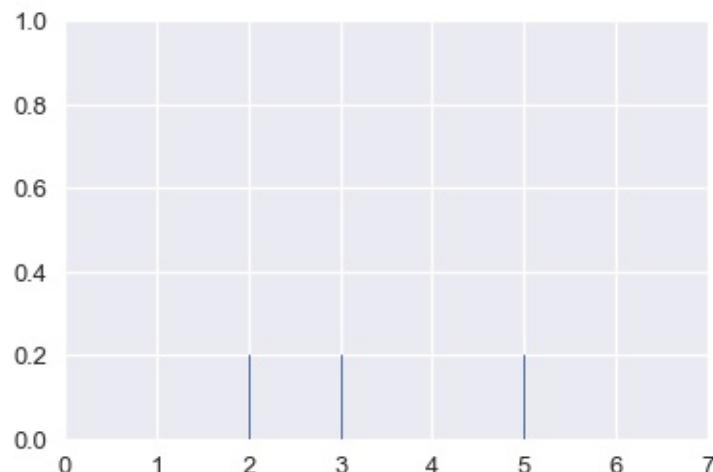
We've also seen that `seaborn` will plot a smooth curve over a histogram by default.

```
sns.distplot(ages)
```

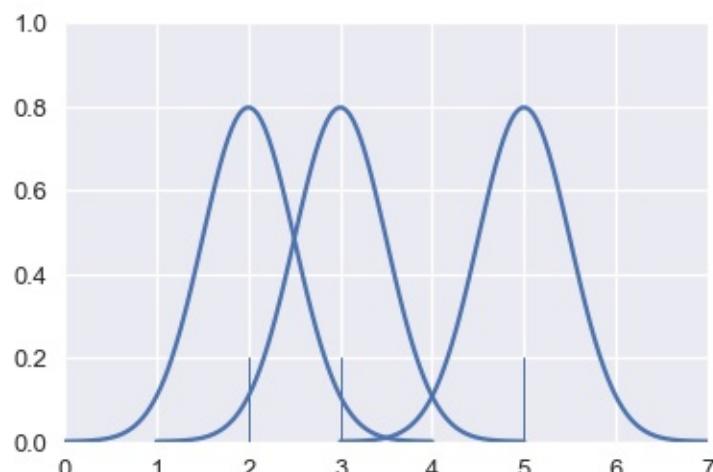
```
<matplotlib.axes._subplots.AxesSubplot at 0x1a23d89780>
```



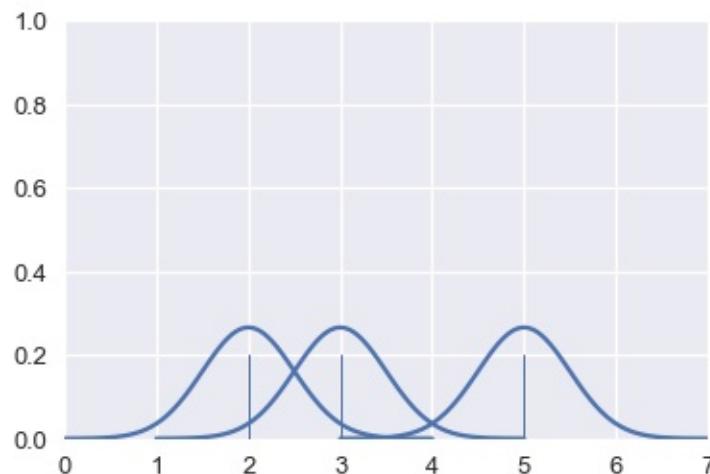
This is another form of smoothing called *kernel density estimation* (KDE). Instead of grouping points together and plotting bars, KDE places a curve on each point and combines the individual curves to create a final estimation of the distribution. Consider the rugplot below that shows three points.



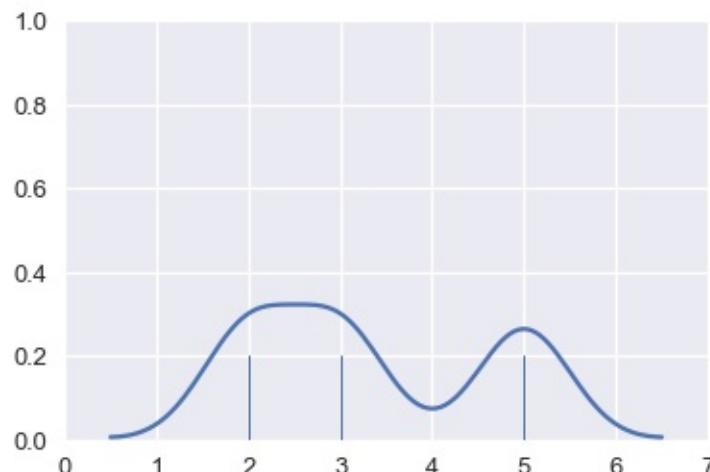
To perform KDE, we place a Gaussian (normal) distribution on each point:



The area under each Gaussian curve is equal to 1. Since we will sum multiple curves together, we scale each curve so that when added together the area under all the curves is equal to 1.

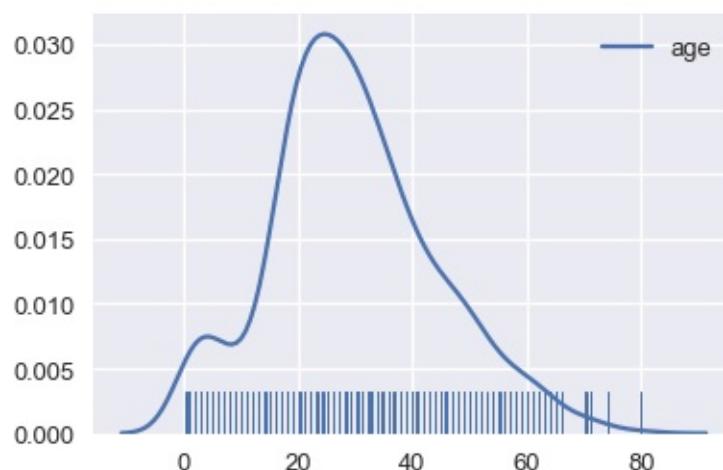


Finally, we add the curves together to create a final smooth estimate for the distribution:



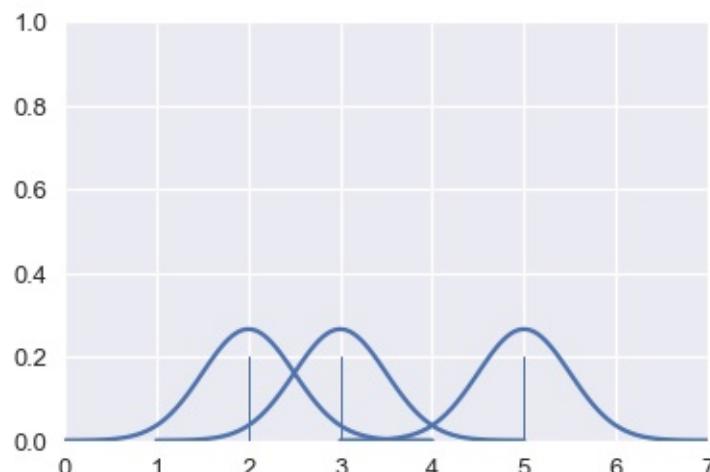
By following this procedure, we can use KDE to smooth many points.

```
# Show the original unsmoothed points  
sns.rugplot(ages, height=0.1)  
  
# Show the smooth estimation of the distribution  
sns.kdeplot(ages);
```

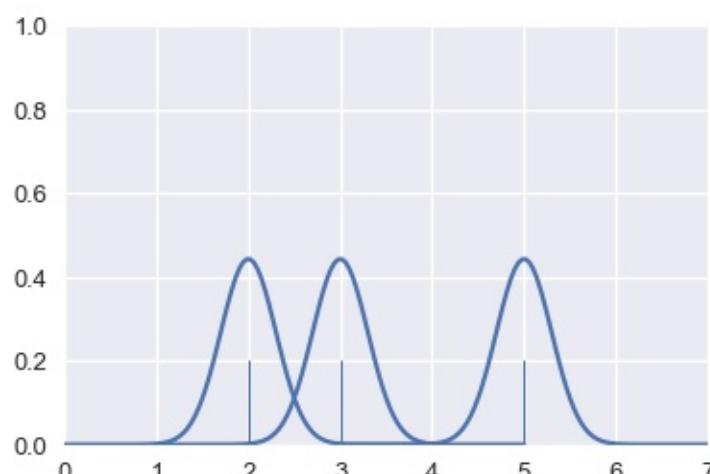


Kernel Density Estimation Details

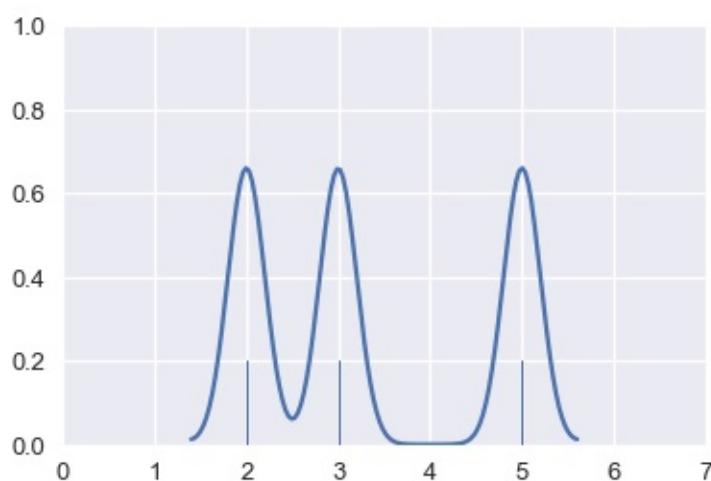
In the previous examples of KDE, we placed a miniature Gaussian curve on each point and added the Gaussians together.



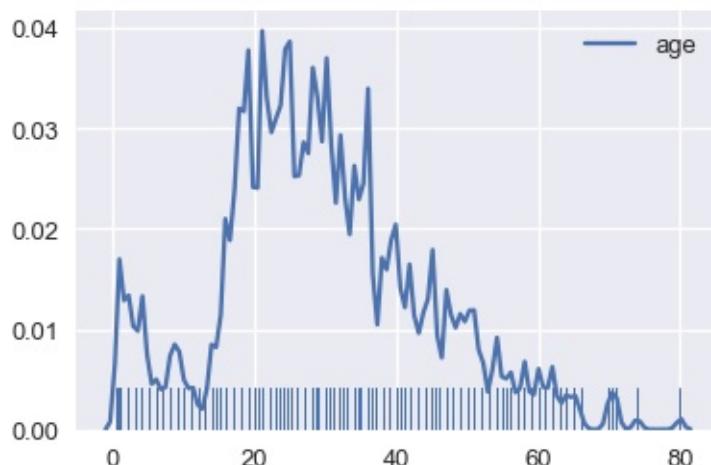
We are free to adjust the width of the Gaussians. For example, we can make each Gaussian narrower. This is called decreasing the *bandwidth* of the kernel estimation.



When we add these narrower Gaussians together, we create a more detailed final estimation.



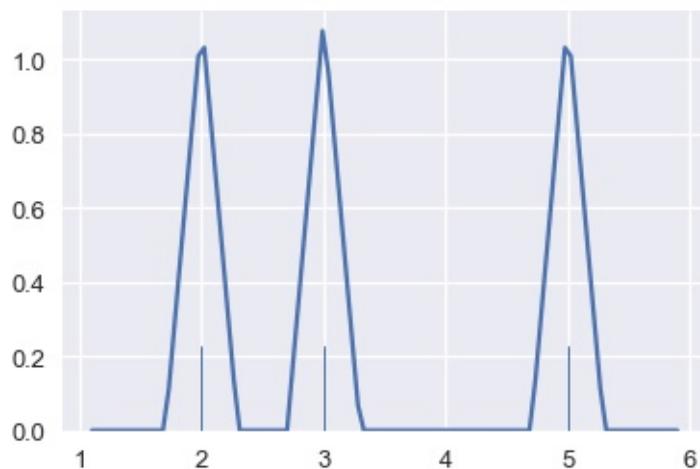
```
# Plot the KDE for Titanic passenger ages using a lower
bandwidth
sns.rugplot(ages, height=0.1)
sns.kdeplot(ages, bw=0.5);
```



Just like adjusting bins for a histogram, we typically adjust the bandwidth until we believe the final plot shows the distribution without distracting the viewer with too much detail.

Although we have placed a Gaussian at each point so far, we can easily select other functions to estimate each point. This is called changing the *kernel* of the kernel density estimation. Previously, we've used a Gaussian kernel. Now, we'll use a triangular kernel which places a pair of stepwise sloped lines at each point:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1a26eb3908>
```



```
# Plot the KDE for Titanic passenger ages using a triangular kernel
sns.rugplot(ages, height=0.1)
sns.kdeplot(ages, kernel='tri');
```



Usually we'll use a Gaussian kernel unless we have a specific reason to use a different kernel.

Smoothing a Scatter Plot

We can also smooth two-dimensional plots when we encounter the problem of overplotting.

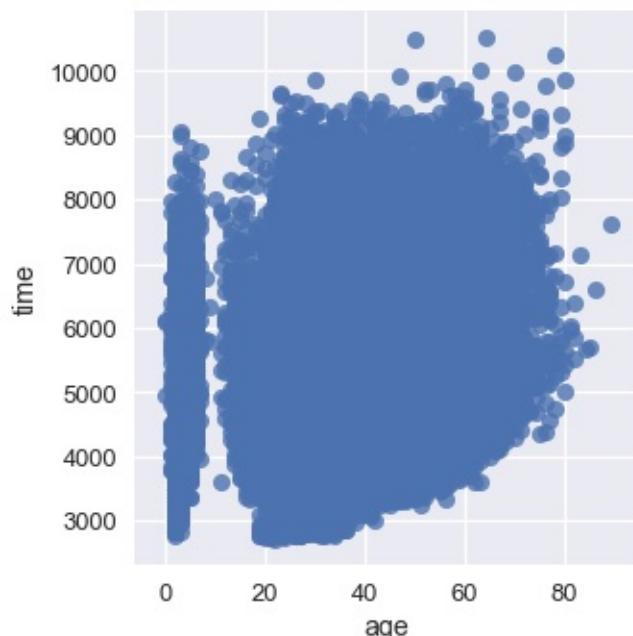
The following example comes from a dataset released by the Cherry Blossom Run, an annual 10-mile run in Washington D.C. Each runner can report their age and their race time; we've plotted all the reported data points in the scatter plot below.

```
runners = pd.read_csv('data/cherryBlossomMen.csv').dropna()
runners
```

	year	place	age	time
0	1999	1	28.0	2819.0
1	1999	2	24.0	2821.0
2	1999	3	27.0	2823.0
...
70066	2012	7190	56.0	8840.0
70067	2012	7191	35.0	8850.0
70069	2012	7193	48.0	9059.0

70045 rows × 4 columns

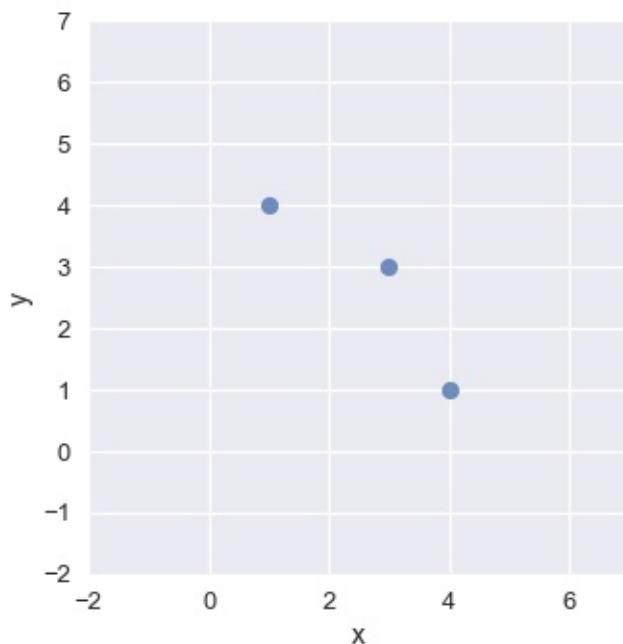
```
sns.lmplot(x='age', y='time', data=runners, fit_reg=False);
```



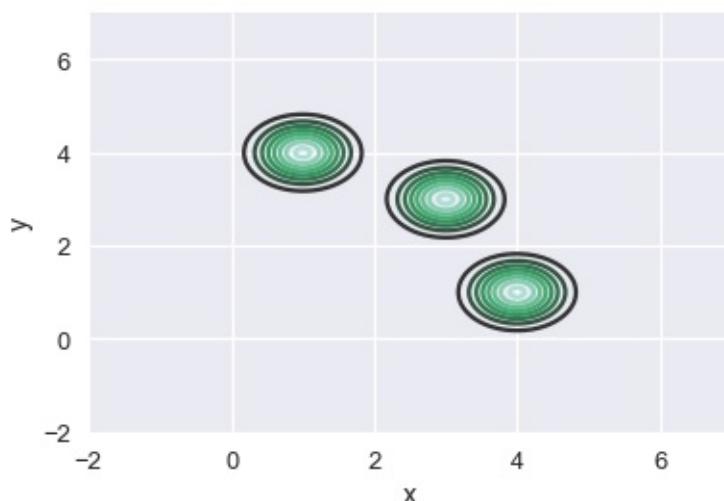
So many points lie on top of each other that it's difficult to see any trend at all!

We can smooth the scatter plot using kernel density estimation in two dimensions. When KDE is applied to a two-dimensional plot, we place a three-dimensional Gaussian at each point. In three dimensions, the Gaussian looks like a mountain pointing out of the page.

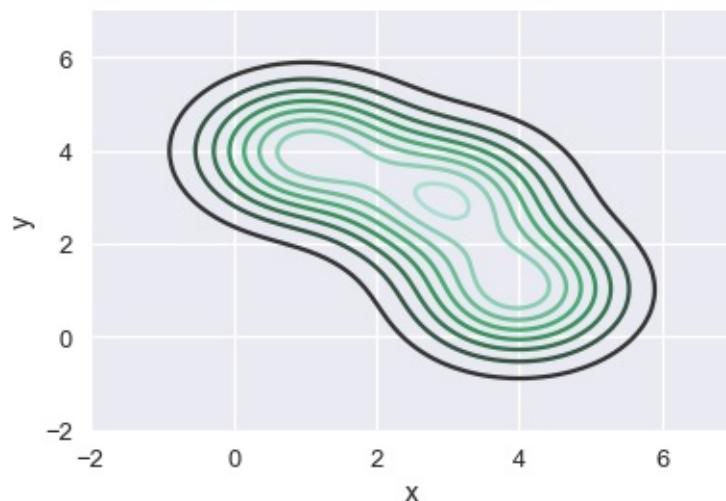
```
# Plot three points
two_d_points = pd.DataFrame({'x': [1, 3, 4], 'y': [4, 3, 1]})
sns.lmplot(x='x', y='y', data=two_d_points, fit_reg=False)
plt.xlim(-2, 7)
plt.ylim(-2, 7);
```



```
# Place a Gaussian at each point and use a contour plot to show
# each one
sns.kdeplot(two_d_points['x'], two_d_points['y'], bw=0.4)
plt.xlim(-2, 7)
plt.ylim(-2, 7);
```

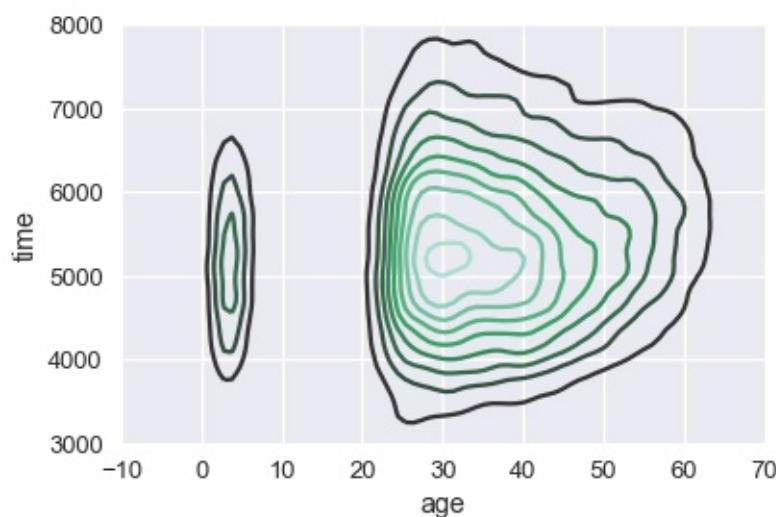


Just like we've previously seen, we scale each Gaussian and add them together to obtain a final contour plot for the scatter plot.



The resulting plot shows the downward sloping trend of the three points. Similarly, we can apply a KDE to smooth out the scatter plot of runner ages and times:

```
sns.kdeplot(runners['age'], runners['time'])
plt.xlim(-10, 70)
plt.ylim(3000, 8000);
```



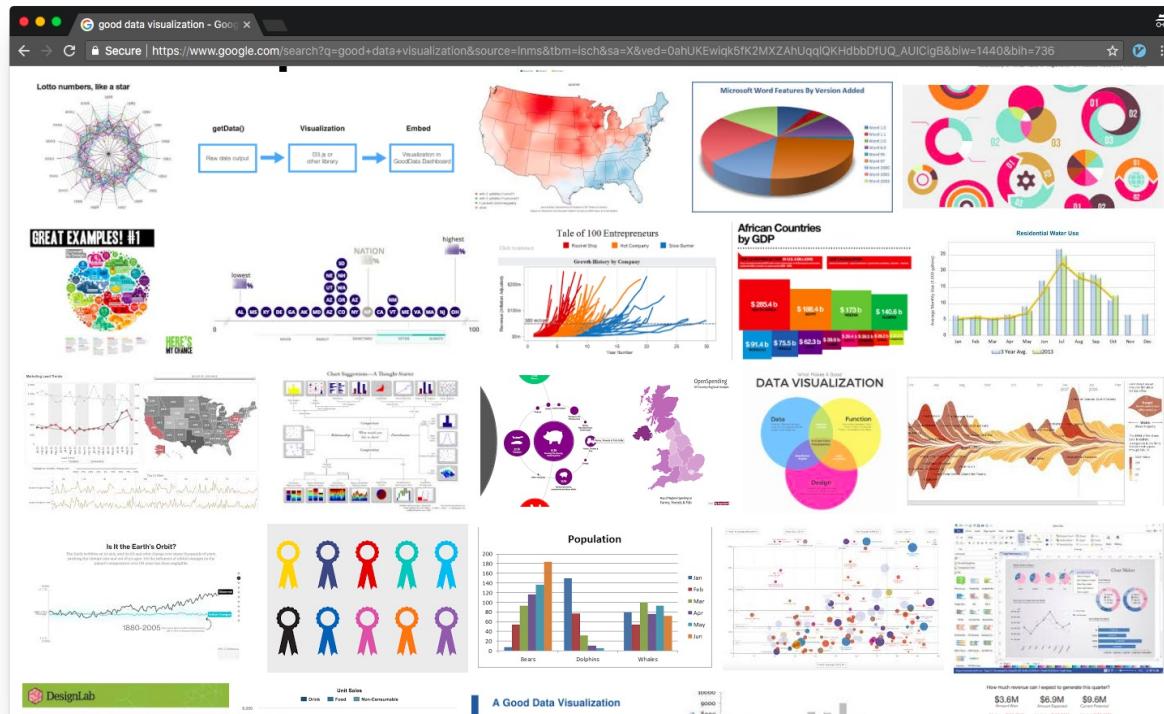
We can see that most of our runners were between 25 and 50 years old, and that most runners took between 4000 and 7000 seconds (roughly between 1 and 2 hours) to finish the race.

We can see more clearly that there is a suspicious group of runners that are between zero and ten years old. We might want to double check that our data for those ages was recorded properly.

We can also see a slight upward trend in the time taken to finish the race as runner age increases.

Philosophy for Data Visualization

Unfortunately, many charts in the world at large do not adhere to our principles for data visualization. Here is a page of charts that appear on the first page of a Google search for "good data visualization" conducted in Spring of 2018. How many poorly made charts can you spot?



Effective data visualization reveals the data. All of our principles for visualization aim to make the data more understandable for the viewer. Data visualization is sometimes used to mislead and misinform. When done properly, however, data visualization is one of our greatest tools for discovering, revealing, and communicating trends and anomalies in our data.

Modeling and Estimation

Essentially, all models are wrong, but some are useful.

— [George Box, Statistician \(1919-2013\)](#)

We have covered question formulation, data cleaning, and exploratory data analysis, the first three steps of the data science lifecycle. We have also seen that EDA often reveals relationships between variables in our dataset. How do we decide whether a relationship is real or spurious? How do we use these relationships to make reliable predictions about the future? To answer these questions we will need the mathematical tools for modeling and estimation.

A model is an **idealized** representation of a system. For example, if we drop a ball off the Leaning Tower of Pisa, we expect the ball to drop to the ground because we have a model of gravity. Our model of gravity also allows us to predict how long it will take the ball to hit the ground using the laws of projectile motion.

This model represents our system but is simplistic—for example, it leaves out the effects of air resistance, the gravitational effects of other celestial bodies, and the buoyancy of air. Because of these unconsidered factors, our model will almost always make incorrect predictions in real life! Still, the simple model of gravity is accurate enough in so many situations that it's widely used and taught today.

In data science, we often use data to create models to estimate how the world will behave in the future. How do we choose a model? How do we decide whether we need a more complicated model? We will explore these questions in this chapter.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Models](#)
 - [A Simple Model](#)
 - [Cost Function Intuition](#)

Models

In the United States, many diners will leave a tip for their waiter or waitress as the diners pay for the meal. Although it is customary to offer 15% of the total bill as tip, perhaps some restaurants have more generous patrons than others.

One particular waiter was so interested in how much tip he could expect to get that he collected a simple random sample of information about the tables he served.

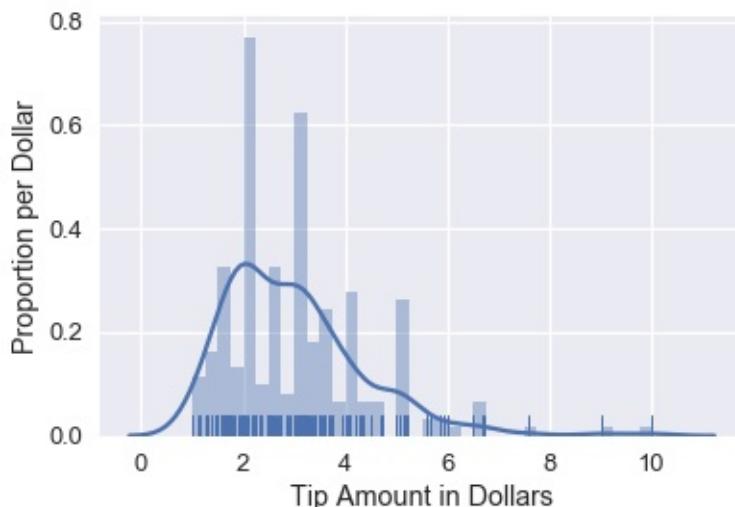
```
tips = sns.load_dataset('tips')
tips
```

	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
...
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

We can plot a histogram of the tip amounts:

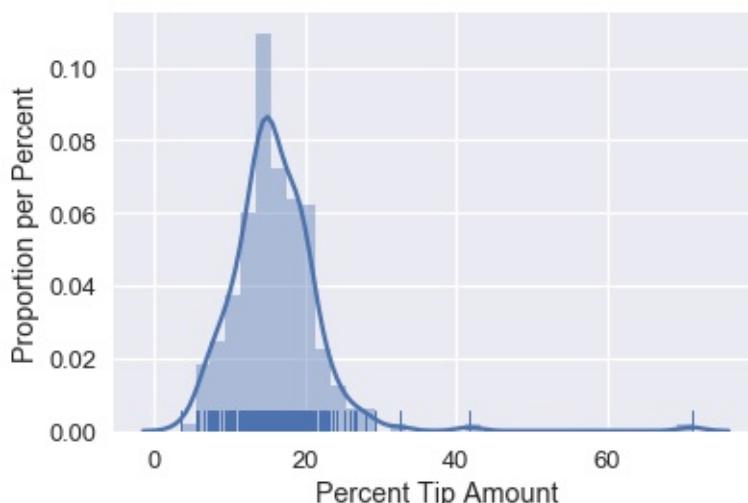
```
sns.distplot(tips['tip'], bins=np.arange(0, 10.1, 0.25),
rug=True)
plt.xlabel('Tip Amount in Dollars')
plt.ylabel('Proportion per Dollar');
```



There are already some interesting patterns in the data. For example, there is a clear mode at \\$2 and most tips seem to be in multiples of \\$0.50.

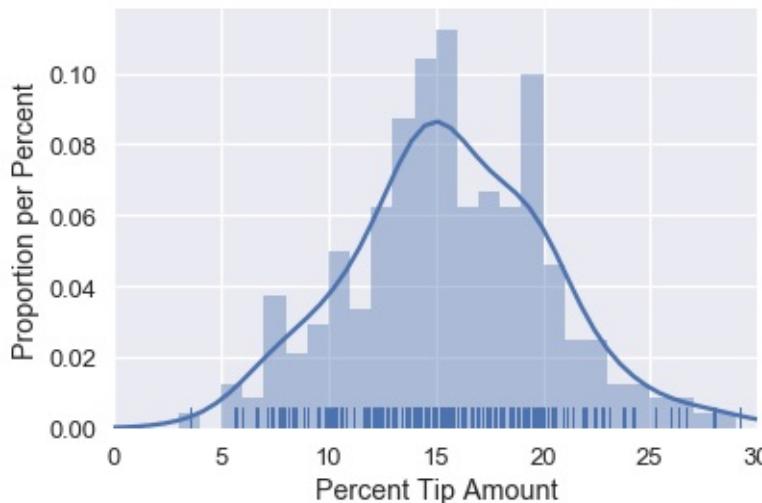
For now, we are most interested in the percent tip: the tip amount divided by the bill amount. We can create a column in our DataFrame for this variable and show its distribution.

```
tips['pcttip'] = tips['tip'] / tips['total_bill'] * 100
sns.distplot(tips['pcttip'], rug=True)
plt.xlabel('Percent Tip Amount')
plt.ylabel('Proportion per Percent');
```



It looks like one table left our waiter a tip of 70%! However, most of the tips percentages are under 30%. Let's zoom into that part of the distribution.

```
sns.distplot(tips['pcttip'], bins=np.arange(30), rug=True)
plt.xlim(0, 30)
plt.xlabel('Percent Tip Amount')
plt.ylabel('Proportion per Percent');
```



We can see the distribution is roughly centered at 15% with another potential mode at 20%. Suppose our waiter is interested in predicting how much percent tip he will get from a given table. To address this question, we can create a model for how much tip the waiter will get.

A Simple Model

The simplest model possible is to completely ignore the data and state that since the convention in the U.S. is to give 15% tip, the waiter will always get 15% tip from his tables. While simple, let's use this model to define some variables that we'll use later on.

This model assumes that there is one true percentage tip that all tables, past and future, will give the waiter. This is the *population parameter* for the percent tip, which we will denote by θ^* .

After making this assumption, our model then says that our guess for θ^* is 15%. We will use θ to represent our estimate.

In mathematical notation, our model states that:

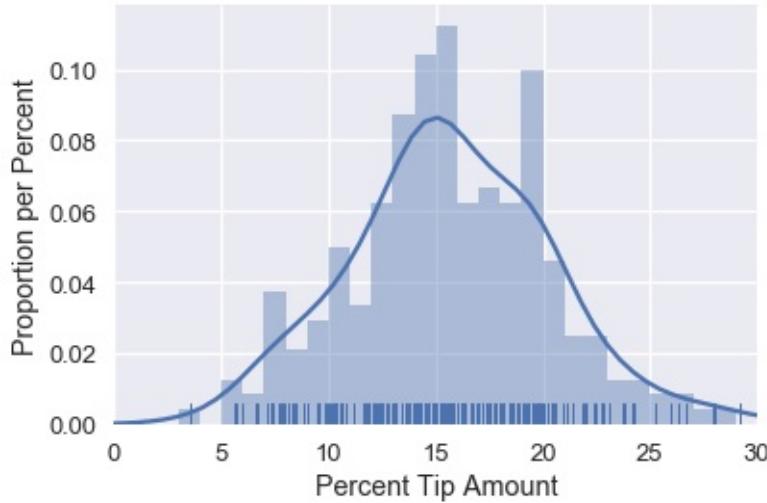
$$\theta = 15$$

This model is obviously wrong—if the model were accurate, every table in our dataset should have given the waiter exactly 15% tip. However, this model will probably make a reasonable guess for most scenarios. In fact, if we had no data on previous tables and their percent tips this model would be a good choice. Still, the question remains: what if 10% is a better estimate? Or 20%? Without data, we are mostly limited to blind guesses.

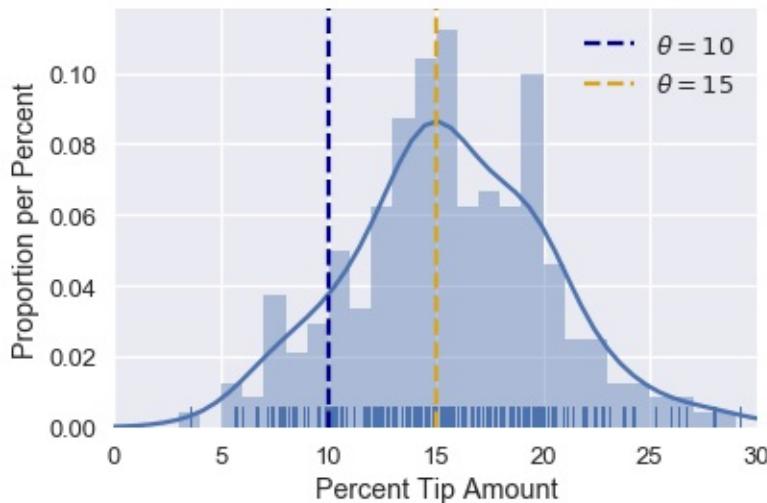
However, since our waiter's dataset was sampled at random from our population of interest, our sample will likely resemble the population. This key insight allows us to use our sample to evaluate the accuracy of our model. In particular, we can use the sample of tip percentages to decide whether 10%, 15%, 20%, or some other percent is the best choice for our estimate θ .

Cost Function Intuition

The distribution of tip percents from our dataset is replicated below for convenience.



Let's suppose we are trying to compare two choices for θ : 10% and 15%. We can mark both of these choices on our distribution:



Intuitively, it looks like choosing $\theta = 15$ makes more sense than $\theta = 10$ given our dataset. Why is this? When we look at the points in our data, we can see that more points fall close to 15 than they do to 10.

Although in this case it is clear that $\theta = 15$ is a better choice than $\theta = 10$, it is not so clear whether $\theta = 15$ is a better choice than $\theta = 16$. By assigning a specific number that measures how "good" a given choice of θ is, we make precise choices between different values of θ . This is a mathematical function that takes in a value of θ and the points in our dataset, outputting a single number that we will use to select the best value of θ that we can.

We call this function a *cost function*.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Cost Functions](#)
 - [Our First Cost Function: Mean Squared Error](#)
 - [A Shorthand](#)
 - [Finding the Best Value of \$\theta\$](#)
 - [The Minimizing Value of the Mean Squared Error](#)
 - [Back to the Original Dataset](#)
 - [Summary](#)

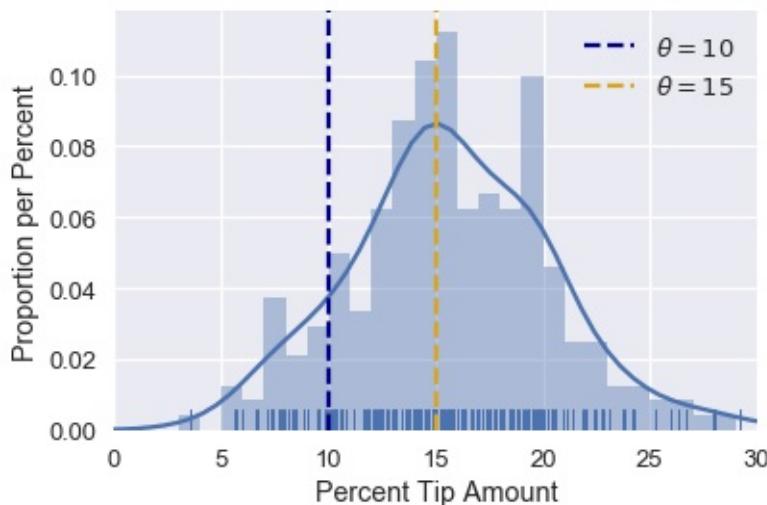
Cost Functions

For now, our model assumes that there is a single population tip percentage θ^* . We are trying to estimate this parameter, and we use the variable θ to denote our estimate. Since our sample of tips is a random sample drawn from the population, we believe that using our sample to create an estimate θ will give a value close to θ^* .

To precisely decide which value of θ is best, we need to define a *cost function*. A cost function is a mathematical function that takes in an estimate θ and the points in our dataset y_1, y_2, \dots, y_n . It outputs a single number that we can use to choose between two different values of θ . In mathematical notation, we want to create the function:

$$L(\theta, y_1, y_2, \dots, y_n) = \dots$$

By convention, the cost function outputs lower values for preferable values of θ and larger values for worse values of θ . In the previous section, we compared $\theta = 10$ and $\theta = 15$.



Since $\theta = 15$ falls closer to most of the points, our cost function should output a small value for $\theta = 15$ and a larger value for $\theta = 10$.

Let's use this intuition to create a cost function.

Our First Cost Function: Mean Squared Error

We would like our choice of θ to fall close to the points in our dataset. Thus, we can define a cost function that outputs a larger value as θ gets further away from the points in the dataset. We start with a simple cost function called the *mean squared error*.

Here's the idea:

1. We select a value of θ .
2. For each value in our dataset, take the squared difference between the value and theta: $(y_i - \theta)^2$. Squaring the difference in a simple way to convert negative differences into positive ones. We want to do this because if our point $y_i = 14$, $\theta = 10$ and $\theta = 18$ are equally far away from the point and are thus equally "bad".
3. To compute the final cost, take the average of each of the individual squared differences.

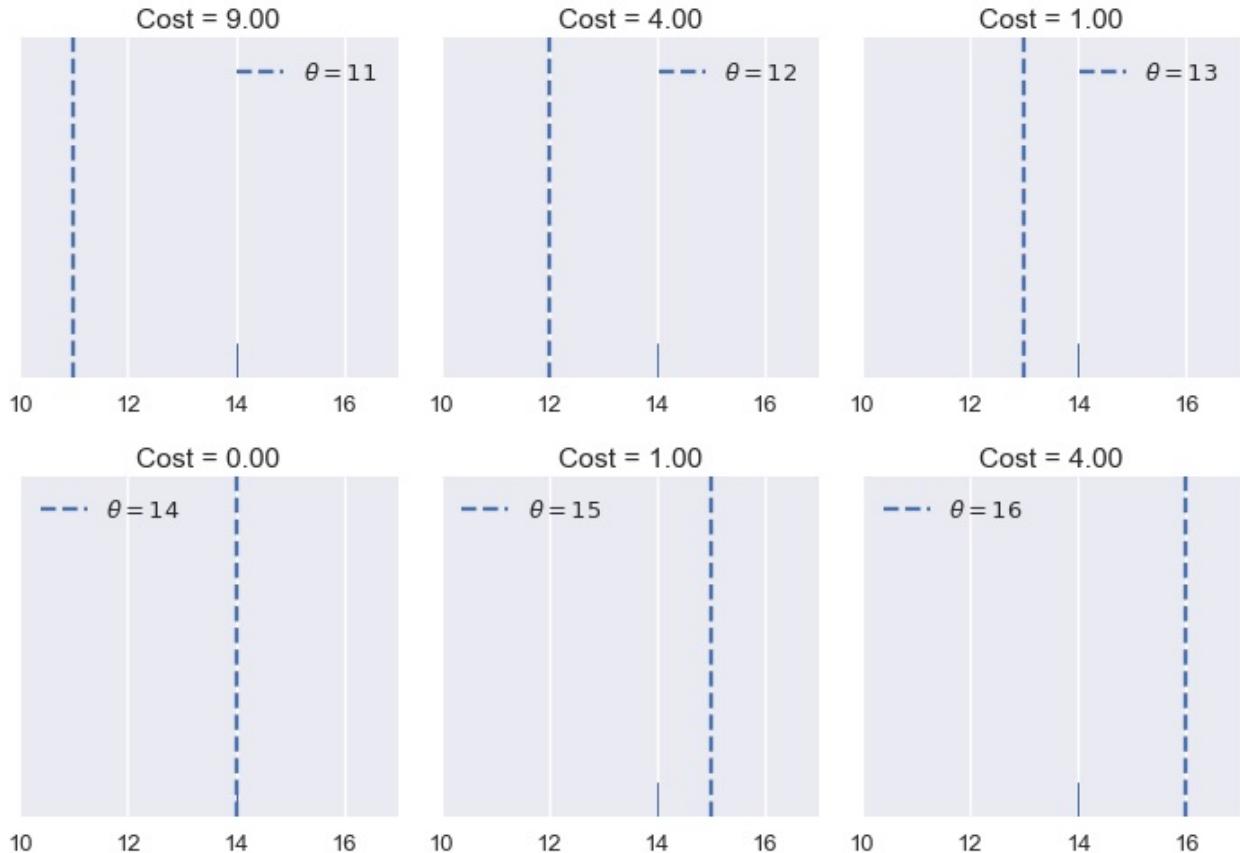
This gives us a final cost function of:

$$\begin{aligned} L(\theta, y_1, y_2, \dots, y_n) &= \text{average} \left\{ (y_1 - \theta)^2, (y_2 - \theta)^2, \dots, (y_n - \theta)^2 \right\} \\ &= \frac{1}{n} \left[(y_1 - \theta)^2 + (y_2 - \theta)^2 + \dots + (y_n - \theta)^2 \right] \end{aligned}$$

Creating a Python function to compute the loss is simple to do:

```
def mse_cost(theta, y_vals):
    return np.mean((y_vals - theta) ** 2)
```

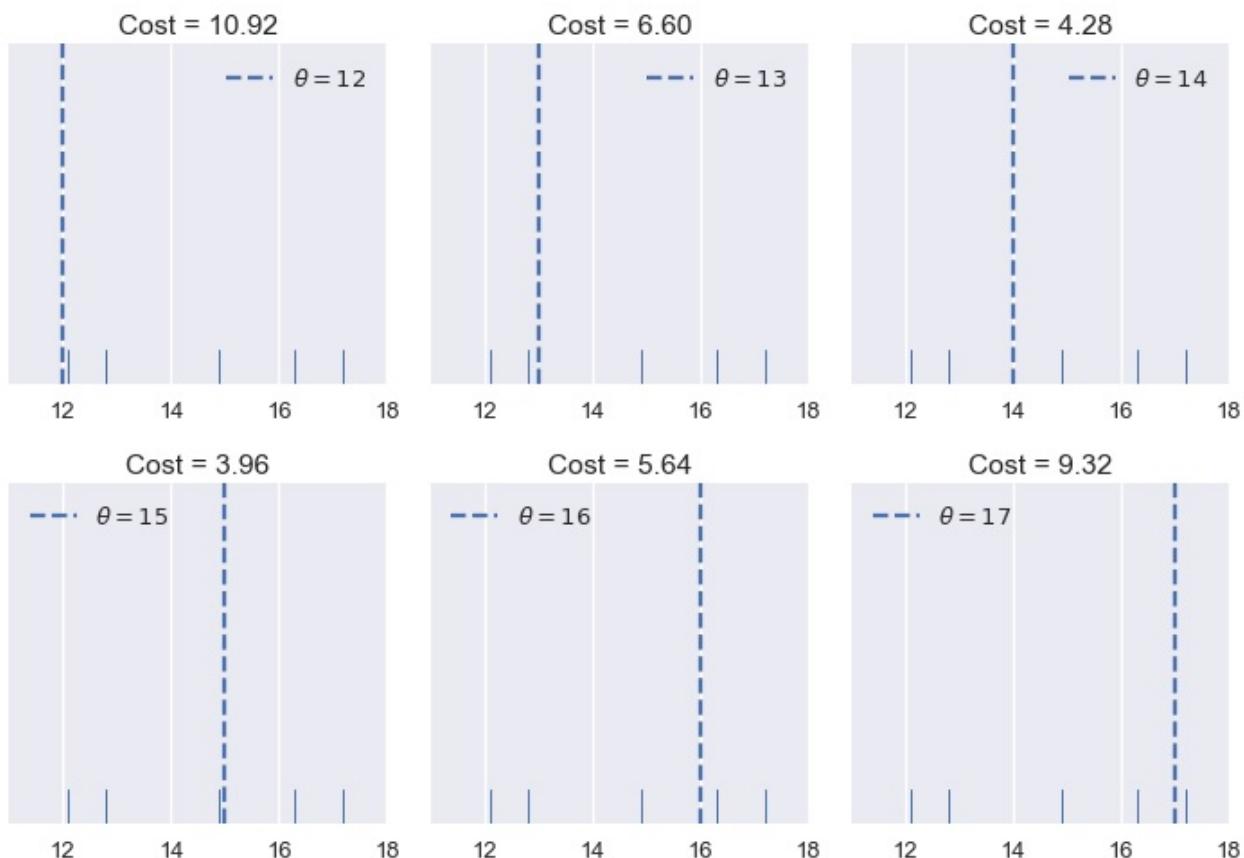
Let's see how this cost function behaves. Suppose we have a dataset only containing one point, $y_1 = 14$. We can try different values of θ and see what the cost function outputs for each value.



You can also interactively try different values of θ below. You should understand why the cost for $\theta = 11$ is many times higher than the cost for $\theta = 13$.

[Show Widget](#)

As we hoped, our cost is larger as θ is further away from our data and is smallest when θ falls exactly onto our data point. Let's now see how our mean squared error behaves when we have five points instead of one. Our data this time are: $\{12.1, 12.8, 14.9, 16.3, 17.2\}$.

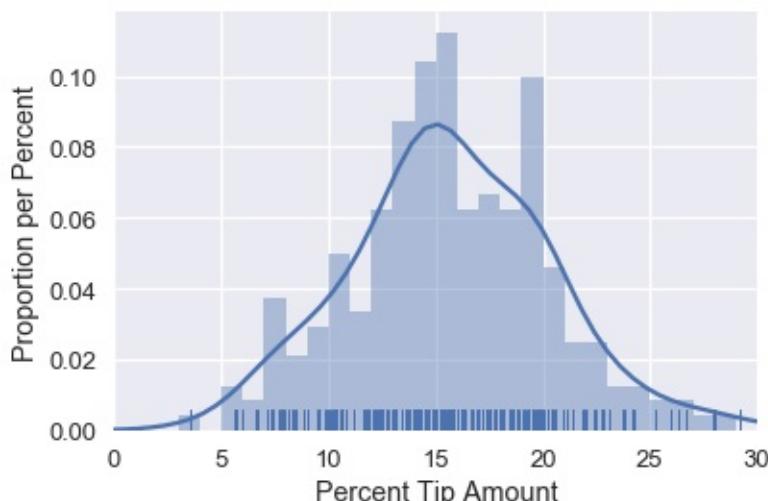


Of the values of θ we tried $\theta = 15$ has the lowest cost. However, a value of θ in between 14 and 15 might have an even lower cost than $\theta = 15$. See if you can find a better value of θ using the interactive plot below.

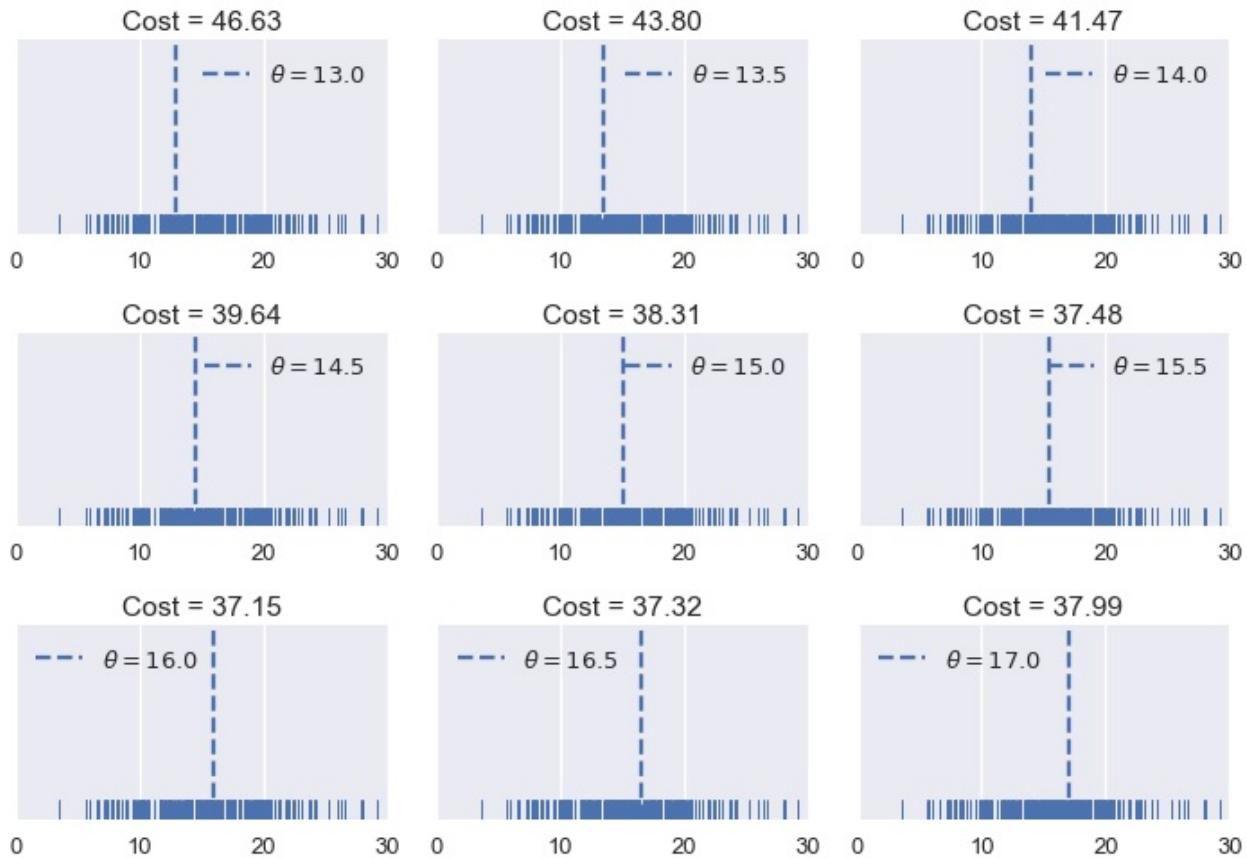
(How might we know for certain whether we've found the best value of θ ? We will tackle this issue soon.)

[Show Widget](#)

The mean squared error cost function seems to be doing its job by penalizing values of θ that are far away from the center of the data. Let's now see what the cost function outputs on the original dataset of tip percents. For reference, the original distribution of tip percents is plotted below:



Let's try some values of θ .



As before, we've created an interactive widget to test different values of θ .

Show Widget

It looks like the best value of θ that we've tried so far is 16.00. This is slightly above our original blind guess of 15% tip. It appears that our waiter gets a bit more tip than we originally thought.

A Shorthand

We have defined our first cost function, the mean squared error cost (MSE). It computes high cost for values of θ that are further away from the center of the data.

Mathematically, this cost function is defined as:

$$\$ \$ \begin{aligned} L(\theta, y_1, y_2, \dots, y_n) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \end{aligned} \$ \$$$

The cost function will compute different costs whenever we change either θ or y_1, y_2, \dots, y_n . We've seen this happen when we tried different values of θ and when we added new data points (changing y_1, y_2, \dots, y_n).

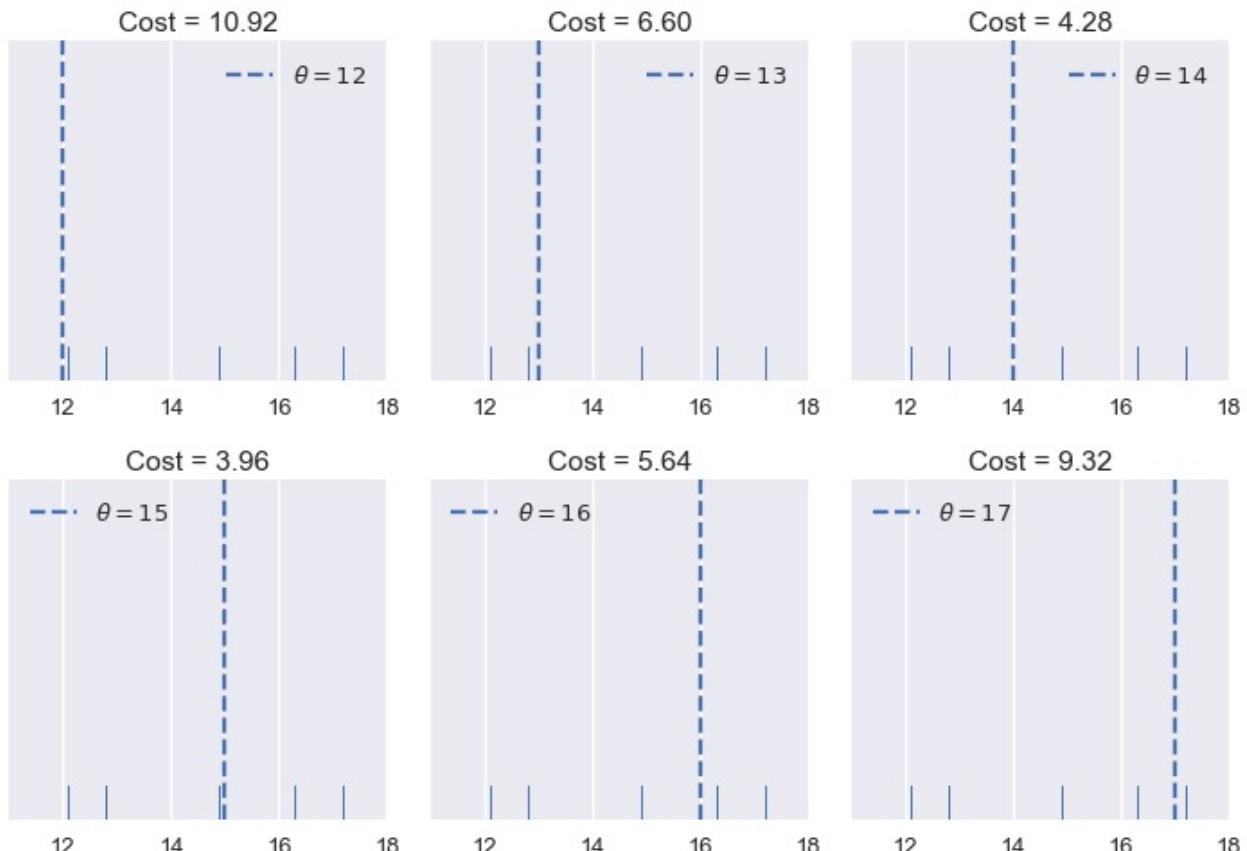
As a shorthand, we can define the vector $y = [y_1, y_2, \dots, y_n]$. Then, we can write our MSE cost as:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \end{aligned} \$ \$$$

Finding the Best Value of θ

So far, we have found the best value of θ by simply trying out a bunch of values and then picking the one with the least cost. Although this method works decently well, we can find a better method by using the properties of our cost function.

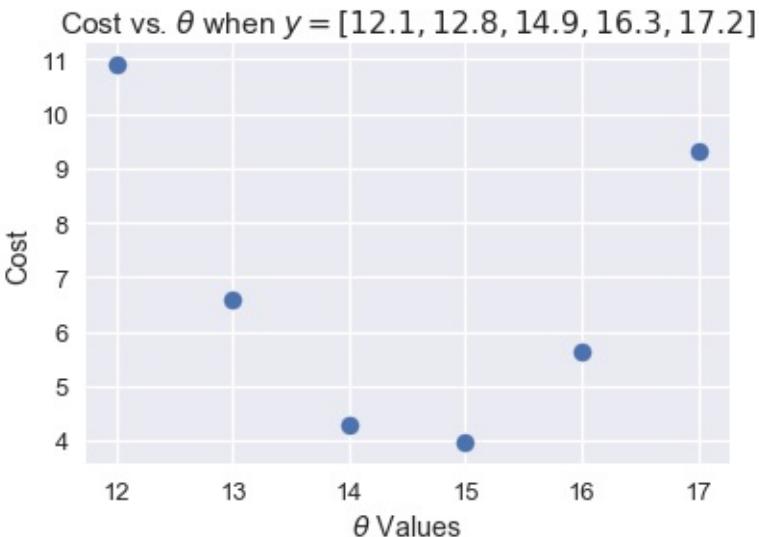
For the time being, let's return to our example with five points: $y = [12.1, 12.8, 14.9, 16.3, 17.2]$.



In the plots above, we've used integer θ values in between 12 and 17. When we change θ , the cost seems to start high (at 10.92), decrease until $\theta = 15$, then increase again. We can see that the cost changes as θ changes, so let's make a plot comparing the cost to θ for each of the six θ s we've tried.

```
thetas = np.array([12, 13, 14, 15, 16, 17])
y_vals = np.array([12.1, 12.8, 14.9, 16.3, 17.2])
costs = [mse_cost(theta, y_vals) for theta in thetas]

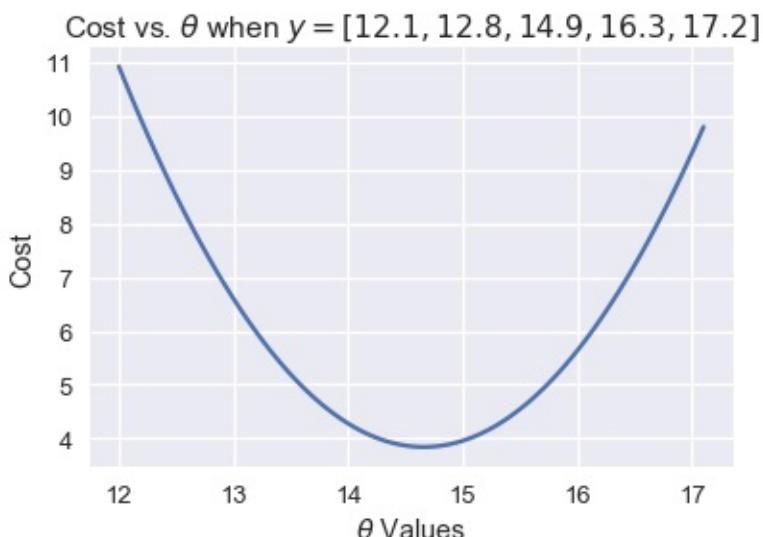
plt.scatter(thetas, costs)
plt.title(r'Cost vs. $\theta$ when $y = [12.1, 12.8, 14.9, 16.3, 17.2]$')
plt.xlabel(r'$\theta$ Values')
plt.ylabel('Cost');
```



The scatter plot shows the downward, then upward trend that we noticed before. We can try more values of θ to see a complete curve that shows how the cost changes as θ changes.

```
thetas = np.arange(12, 17.1, 0.05)
y_vals = np.array([12.1, 12.8, 14.9, 16.3, 17.2])
costs = [mse_cost(theta, y_vals) for theta in thetas]

plt.plot(thetas, costs)
plt.title(r'Cost vs. $\theta$ when $y = [ 12.1, 12.8, 14.9, 16.3, 17.2 ] $')
plt.xlabel(r'$\theta$ Values')
plt.ylabel('Cost');
```



The plot above shows that in fact, $\theta = 15$ was not the best choice; a θ of around 14.7 would have gotten a lower cost! We can use calculus to find that minimizing value of θ exactly. First, we start with our cost function:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \end{aligned} \$ \$$$

And plug in our points $y = [12.1, 12.8, 14.9, 16.3, 17.2]$:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{5} ((12.1 - \theta)^2 + (12.8 - \theta)^2 + (14.9 - \theta)^2 + (16.3 - \theta)^2 + (17.2 - \theta)^2) \end{aligned} \$ \$$$

To find the value of θ that minimizes this function, we compute the derivative with respect to θ :

$$\$ \$ \begin{aligned} \frac{\partial}{\partial \theta} L(\theta, y) &= \frac{1}{5} (-2(12.1 - \theta) - 2(12.8 - \theta) - 2(14.9 - \theta) - 2(16.3 - \theta) - 2(17.2 - \theta)) \\ &= -\frac{2}{5} ((12.1 - \theta) + (12.8 - \theta) + (14.9 - \theta) + (16.3 - \theta) + (17.2 - \theta)) \end{aligned} \$ \$$$

Then, we find the value of θ where the derivative is zero:

$$\$ \$ \begin{aligned} 0 &= -\frac{2}{5} ((12.1 - \theta) + (12.8 - \theta) + (14.9 - \theta) + (16.3 - \theta) + (17.2 - \theta)) \\ 0 &= (12.1 - \theta) + (12.8 - \theta) + (14.9 - \theta) + (16.3 - \theta) + (17.2 - \theta) \\ 5\theta &= 12.1 + 12.8 + 14.9 + 16.3 + 17.2 \\ \theta &= \frac{12.1 + 12.8 + 14.9 + 16.3 + 17.2}{5} \end{aligned} \$ \$$$

As expected, the value of θ that minimizes the cost is between 14 and 15.

Turn your attention to the second-to-last line of the simplification above:

$$\$ \$ \theta = \frac{12.1 + 12.8 + 14.9 + 16.3 + 17.2}{5} \$ \$$$

Notice that this is a familiar expression: it is the average of the five data points. Could this be a pattern for all values of y ?

The Minimizing Value of the Mean Squared Error

We have seen that different values of θ produce different costs when using the mean squared error cost function. The arithmetic above hints that the value of θ that minimizes the cost is the mean of all of the data points. To confirm this, we turn back to the definition of our cost function. Instead of plugging in points, we take the derivative with respect to θ of the cost function as-is:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \\ \frac{\partial}{\partial \theta} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n -2(y_i - \theta) \\ &= -\frac{2}{n} \sum_{i=1}^n (y_i - \theta) \end{aligned} \$ \$$$

Notice that we have left the variables y_i untouched. We are no longer working with the previous example dataset of five points; this equation can be used with any dataset with any number of points.

Now, we set the derivative equal to zero and solve for θ to find the minimizing value of θ as before:

$$\begin{aligned} -\frac{2}{n} \sum_{i=1}^n (y_i - \theta) &= 0 \\ \sum_{i=1}^n (y_i - \theta) &= 0 \\ \sum_{i=1}^n y_i - n\theta &= 0 \\ \sum_{i=1}^n y_i &\approx \theta + \dots + y_n \\ \theta &\approx \frac{y_1 + \dots + y_n}{n} \end{aligned}$$

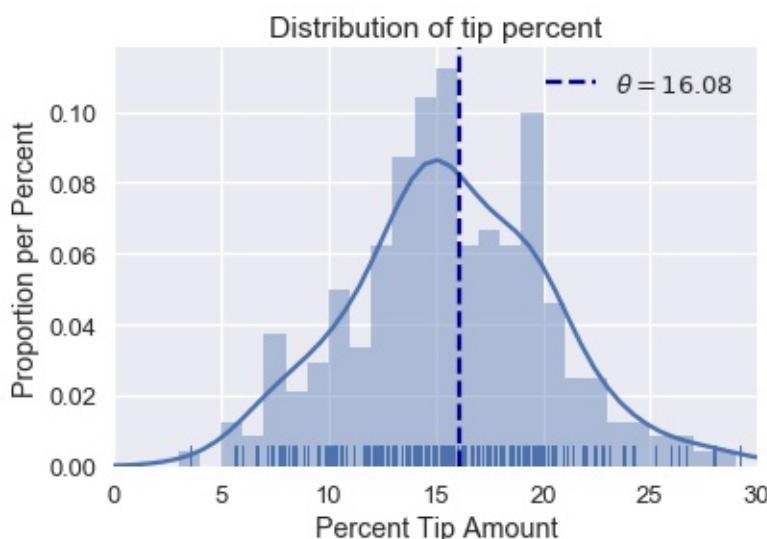
Lo and behold, we see that there is a single value of θ that gives the least MSE no matter what the dataset is. For the mean squared error, we can set θ equal to the mean of the dataset and be confident knowing that we have minimized the cost.

Back to the Original Dataset

We no longer have to test out different values of θ as we did before. We can compute the mean tip percentage in one go:

```
np.mean(tips['pcttip'])
```

```
16.080258172250463
```



Summary

First, we restricted our model to only make a single number as its prediction for all tables. Next, we assume that the waiter's dataset of tips is similar to the population distribution of tip percentages. If this assumption holds, predicting 16.08% will give us the most accurate predictions that we can given our data.

To be more precise, we say that the model is accurate if it minimizes the squared difference between the predictions and the actual values.

Although our model is simple, it illustrates concepts that we'll see over and over again. Future chapters will introduce complicated models. Still, we will discuss each model's assumptions, define cost functions, and find the model that minimizes the cost. It is very helpful to understand this process for simple models before attempting to understand complex ones.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Absolute and Huber Cost](#)
 - [Comparing MSE and Mean Absolute Cost](#)
 - [Outliers](#)
 - [Minimizing the Mean Absolute Loss](#)
 - [MSE and Mean Absolute Cost Comparison](#)
 - [The Huber Cost](#)

Absolute and Huber Cost

Previously, we said that our model is accurate if it minimizes the squared difference between the predictions and the actual values. We used the mean squared error (MSE) cost to capture this measure of accuracy:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \end{aligned} \$ \$$$

We used a simple model that always predicts the same number:

$$\$ \$ \theta = C \$ \$$$

Where C is some constant. When we use this constant model and the MSE cost function, we found that C will always be the mean of the data points. When applied to the tips dataset, we found that the constant model should predict 16.08% since 16.08% is the mean of the tip percents.

Now, we will keep our model the same but switch to a different cost function: the mean absolute cost. Instead taking the squared difference for each point and our prediction, this cost function takes the absolute difference:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n |y_i - \theta| \end{aligned} \$ \$$$

Brief Aside: Cost Functions vs. Loss Functions

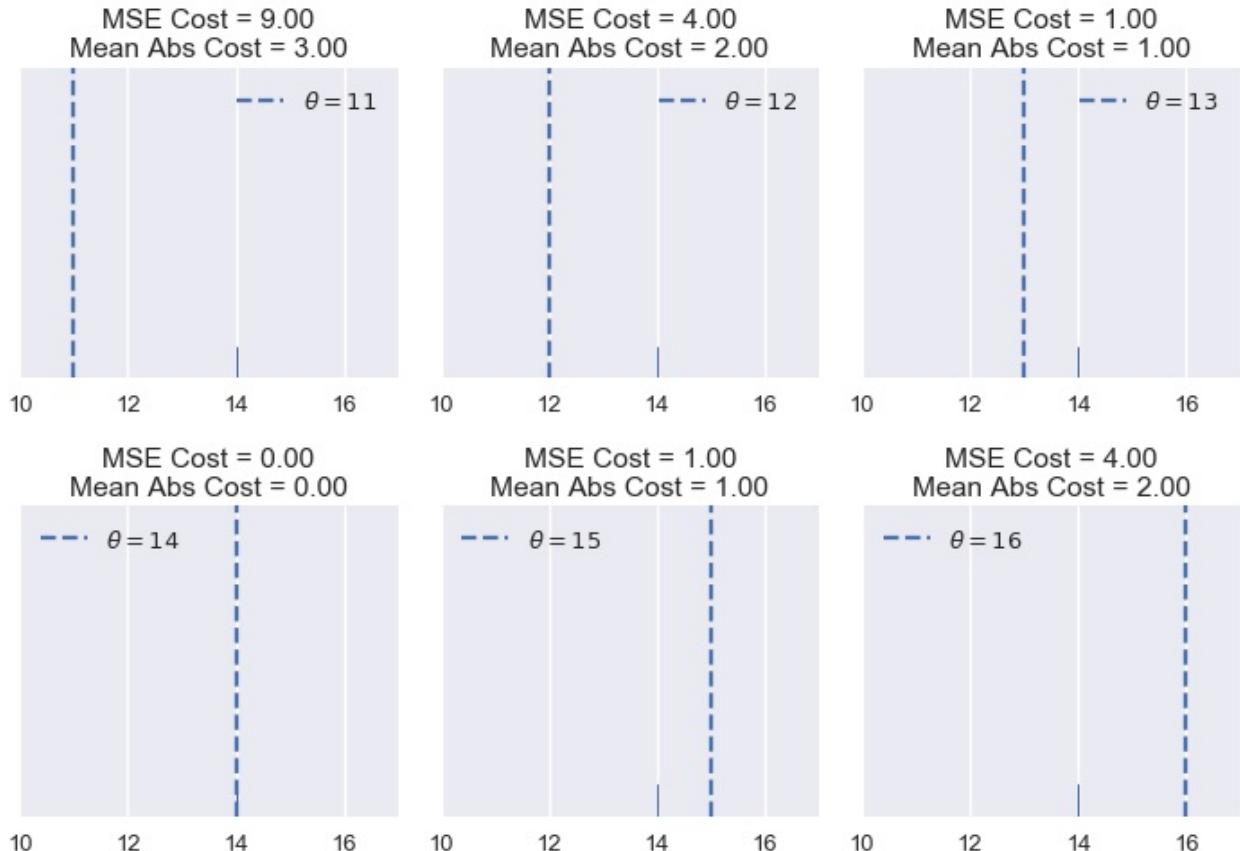
We typically use the term *cost function* to describe a function that takes in parameters of our model (e.g. θ) and the entire set of data that we are interested in (e.g. y). This function generates a single value that describes how well the model does on a given set of points.

We typically use the term *loss function* to describe a function that takes in parameters of our model (e.g. θ) and one point from the data that we are interested in (e.g. y_i). For example, $(y_i - \theta)^2$ is the squared loss, and $|y_i - \theta|$ is the absolute

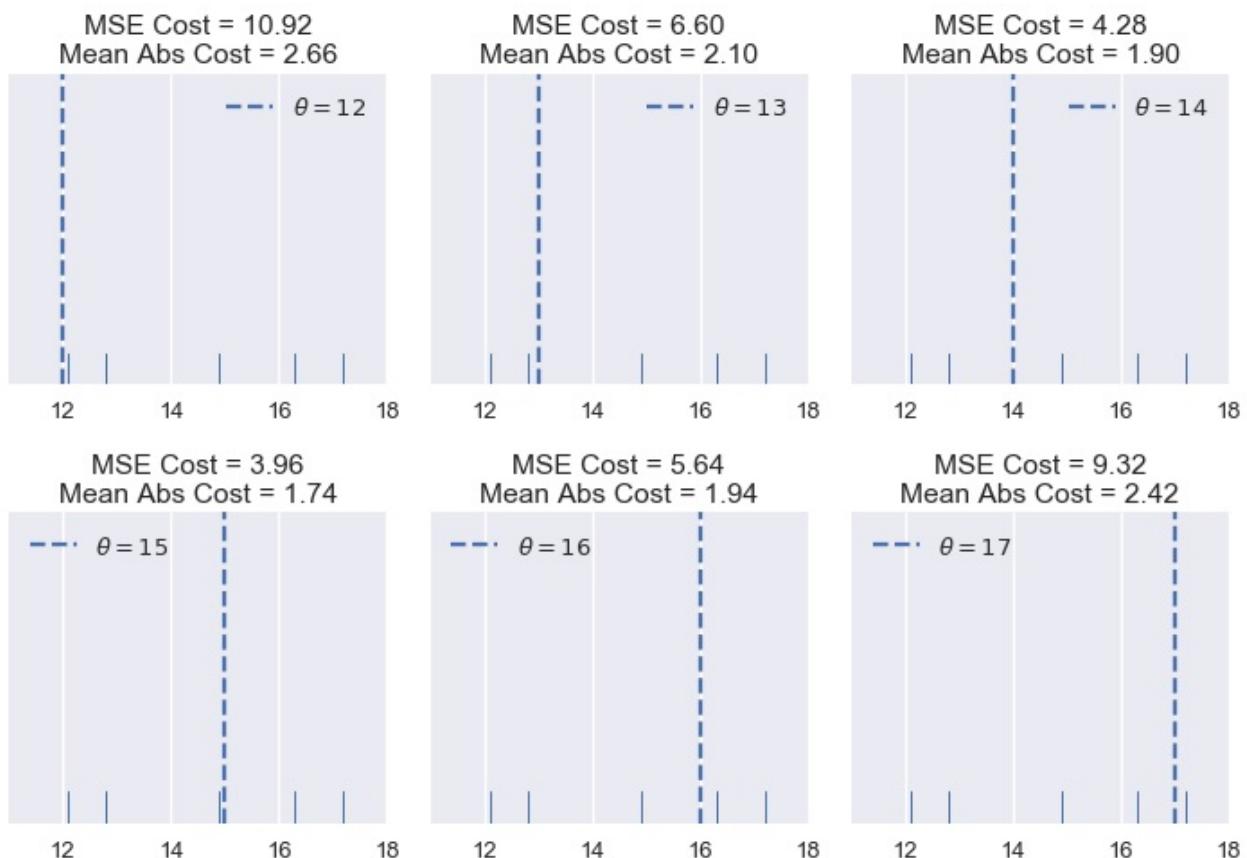
loss. Notice that there is no summation in these example loss functions. You can think of a cost function as the combination of multiple loss function values.

Comparing MSE and Mean Absolute Cost

To get a better sense of how the MSE cost and mean absolute cost compare, let's compare their costs on different datasets. First, we'll use our dataset of one point: $y = [14]$.

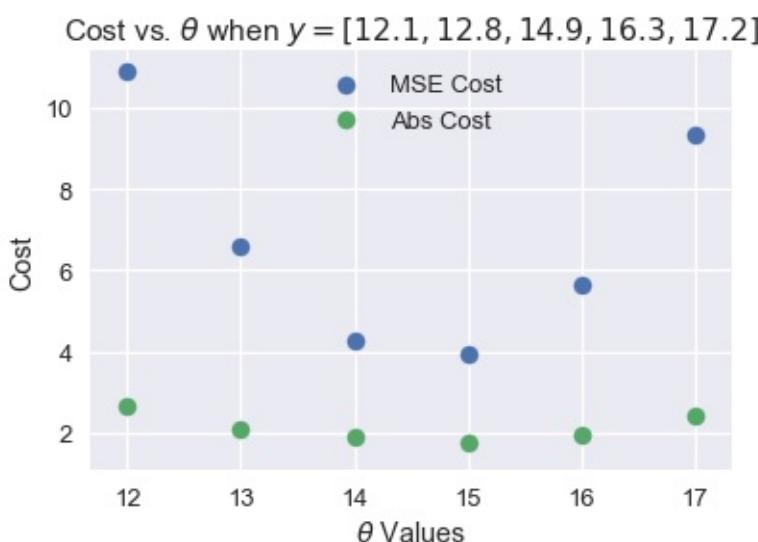


We see that the MSE cost is usually higher than the mean absolute cost since the error is squared. Let's see what happens when have five points: $y = [12.1, 12.8, 14.9, 16.3, 17.2]$

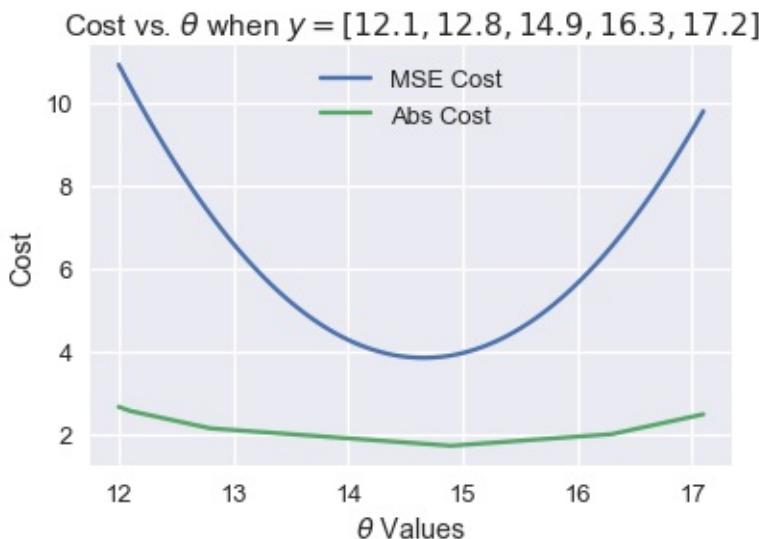


Remember that the actual cost values themselves are not very interesting to us; they are only useful for comparing different values of θ . Once we choose a cost function, we will look for the θ that produces the least cost. Thus, we are interested in whether the cost functions are minimized by different values of θ .

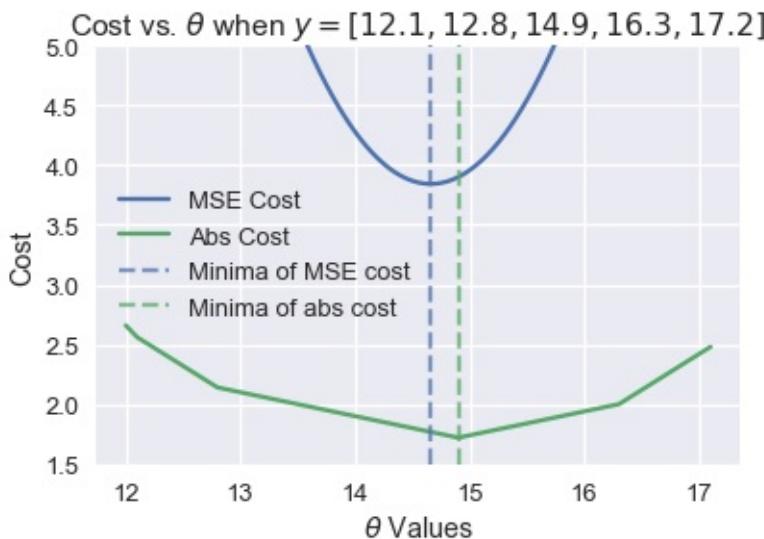
So far, the two cost functions seem to agree on the θ values that produce the least cost for the values that we've tried. If we look a bit closer, however, we will start to see some differences. We first take the costs and plot them against θ for each of the six θ values we tried.



Then, we compute more values of θ so that the curve is smooth:



Then, we zoom into the region between 1.5 and 5 on the y-axis to see the difference in minima more clearly. We've marked the minima with dotted lines.



We've found empirically that the MSE cost and mean absolute cost can be minimized by different values of θ for the same dataset. However, we don't have a good sense of when they will differ and more importantly, why they differ.

Outliers¶

One difference that we can see in the plots of cost vs. θ above lies in the shape of the cost curves. Plotting the MSE cost results in a parabolic curve resulting from the squared term in the cost function.

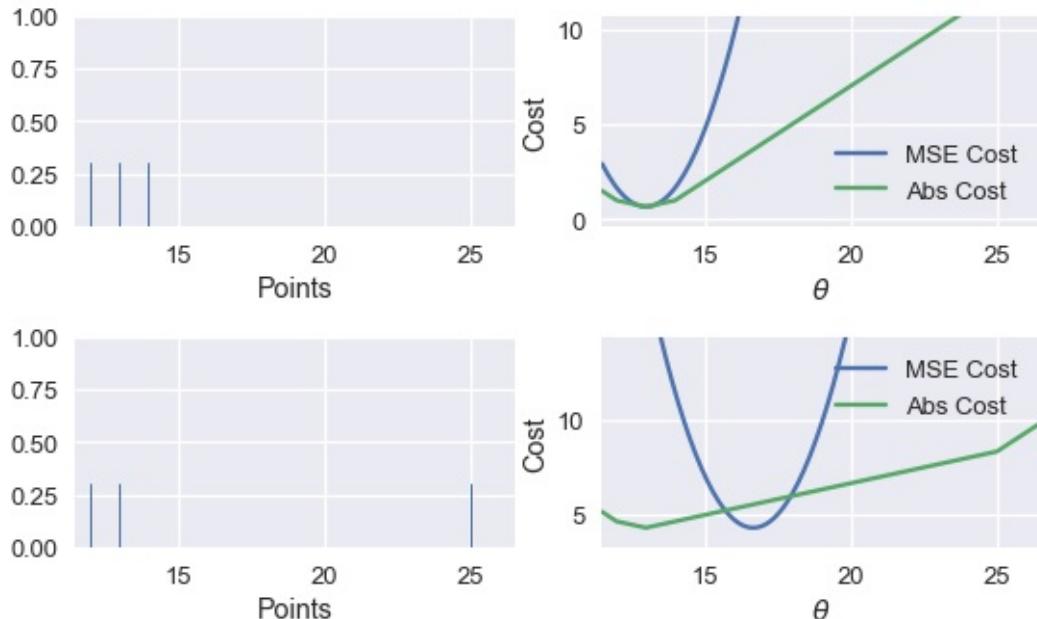
Plotting the mean absolute cost, on the other hand, results in what looks like a connected series of lines. This makes sense when we consider that the absolute value function is linear, so taking the average of many absolute value functions should produce a semi-linear function.

Since the MSE cost has a squared error term, it will be more sensitive to outliers. If $\theta = 10$ and a point lies at 110, that point's error term for MSE will be $(10 - 110)^2 = 10000$ whereas in the mean absolute cost, that point's error term will be $|10 - 110| = 100$. We can illustrate this by taking a set of three points $y = [12, 13, 14]$ and plotting the cost vs. θ curves for MSE and mean absolute loss.

Use the slider below to move the third point further away from the rest of the data and observe what happens to the cost curves. (We've scaled the curves to keep both in view since the MSE cost has larger values than the mean absolute cost.)

Show Widget

We've shown the curves for $y_3 = 14$ and $y_3 = 25$ below.



As we move the point further away from the rest of the data, the MSE cost curve moves with it. When $y_3 = 14$, both MSE and mean absolute cost will be minimized by the same value of θ at $\theta = 13$. However, when $y_3 = 25$, the MSE curve will tell us that the best θ is around 16.7 while the mean absolute cost will still say that $\theta = 13$ is best.

Minimizing the Mean Absolute Loss

Now that we have a qualitative sense of how the MSE and mean absolute cost differ, we can minimize the mean absolute cost to make this difference more precise. As before, we will take the derivative of the cost function with respect to θ and set it equal to zero.

This time, however, we have to deal with the fact that the absolute function is not always differentiable. When $x > 0$, $\frac{\partial}{\partial x} |x| = 1$. When $x < 0$, $\frac{\partial}{\partial x} |x| = -1$. Although $|x|$ is not technically differentiable at $x = 0$, we will set $\frac{\partial}{\partial x} |x| = 0$ so that the equations are easier to work with.

Recall that the equation for the mean absolute loss is:

$$\$ \$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n |y_i - \theta| \\ &= \frac{1}{n} \left(\sum_{y_i < \theta} |y_i - \theta| + \sum_{y_i = \theta} |y_i - \theta| + \sum_{y_i > \theta} |y_i - \theta| \right) \end{aligned} \$ \$$$

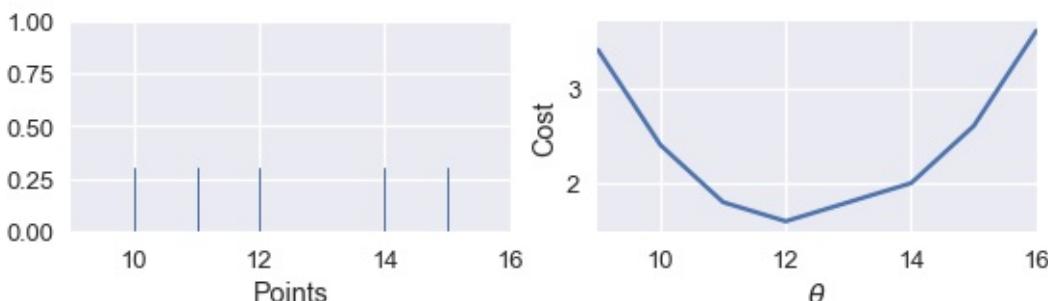
In the line above, we've split up the summation into three separate summations: one that has one term for each $y_i < \theta$, one for $y_i = \theta$, and one for $y_i > \theta$. Why make the summation seemingly more complicated? If we know that $y_i < \theta$ we also know that $|y_i - \theta| < 0$ and thus $\frac{\partial}{\partial \theta} |y_i - \theta| = -1$ from before. A similar logic holds for each term above to make taking the derivative much easier.

Now, we take the derivative with respect to θ and set it equal to zero:

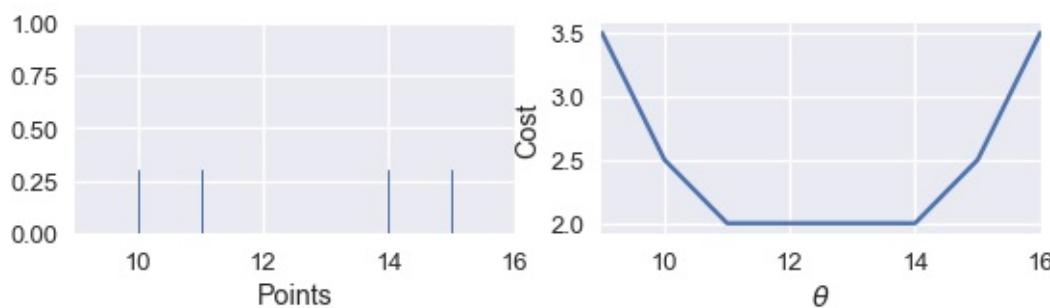
$$\$ \$ \begin{aligned} \frac{1}{n} \left(\sum_{y_i < \theta} (-1) + \sum_{y_i = \theta} (0) + \sum_{y_i > \theta} (1) \right) &= 0 \\ \sum_{y_i < \theta} (-1) + \sum_{y_i > \theta} (1) &= 0 \\ -\sum_{y_i < \theta} (1) + \sum_{y_i > \theta} (1) &= 0 \\ \sum_{y_i < \theta} (1) &= \sum_{y_i > \theta} (1) \end{aligned} \$ \$$$

What does the result above mean? On the left hand side, we have one term for each data point less than θ . On the right, we have one for each data point greater than θ . Then, in order to satisfy the equation we need to pick a value for θ that has the same number of smaller and larger points. This is the definition for the *median* of a set of numbers. Thus, the minimizing value of θ for the mean absolute cost is $\theta = \text{median}(y)$.

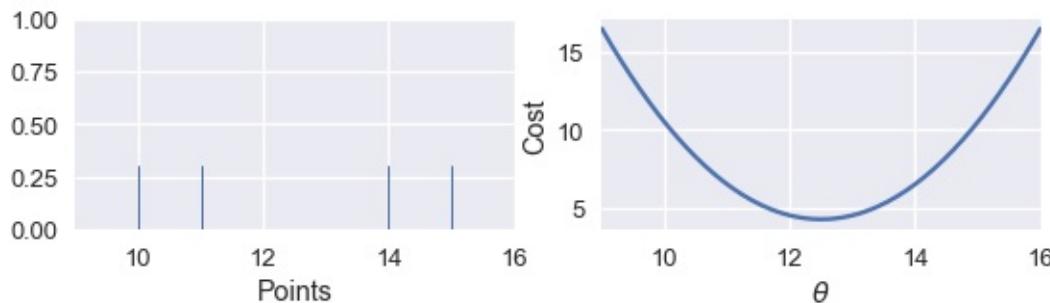
When we have an odd number of points, the median is simply the middle point when the points are arranged in sorted order. We can see that in the example below with five points, the cost is minimized when θ lies at the median:



However, when we have an even number of points, the cost is minimized when θ is any value in between the two central points.



This is not the case when we use the MSE cost:



MSE and Mean Absolute Cost Comparison

Our investigation and the derivation above show that the MSE is easier to differentiate but is more sensitive to outliers than the mean absolute cost. The minimizing θ for MSE is the mean of the data points, and the minimizing θ for the mean absolute cost is the median of the data points. Notice that the median is robust to outliers while the mean is not! This phenomenon arises from our construction of the two cost functions.

We have also seen that the MSE cost will be minimized by a unique value of θ , whereas the mean absolute value can be minimized by multiple values of θ when there are an even number of data points.

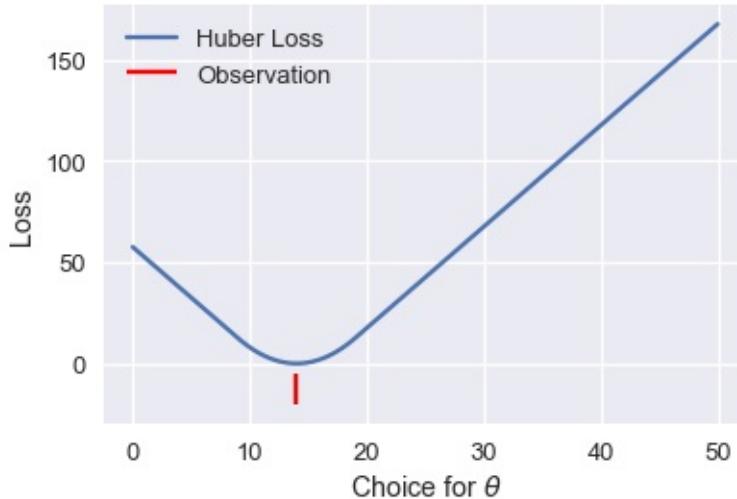
In the examples so far, the ability to differentiate the cost function isn't that useful since we know the exact minimizing value of θ in both cases. However, the ability to differentiate the cost function becomes very important once we start using complicated models. For complicated models, we will not be able to differentiate the cost function by hand and will need a computer to minimize the cost function for us. We will return to this issue when we cover gradient descent and numerical optimization.

The Huber Cost

A third loss function called the Huber loss combines both the MSE and mean absolute loss to create a loss function that is differentiable *and* robust to outliers. The Huber loss accomplishes this by behaving like the MSE loss function at values close to θ and switching to the absolute loss for values far from θ .

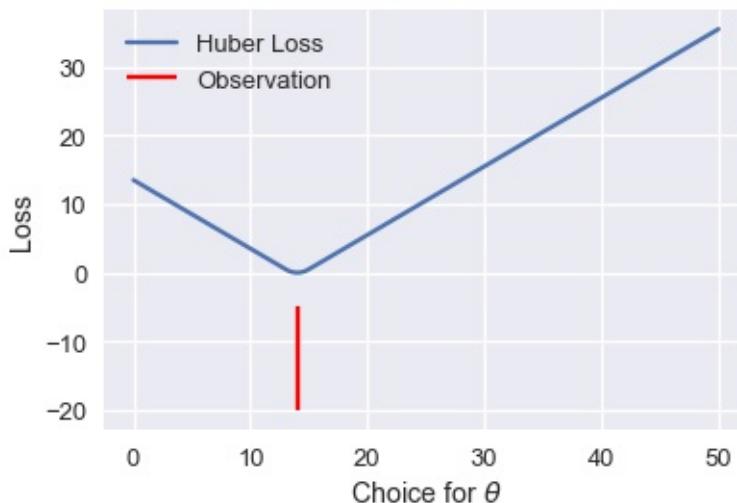
As usual, we create a cost function by taking the mean of the Huber losses for each point in our dataset.

Let's see what the Huber cost function outputs for a dataset of $y = [14]$ as we vary θ :

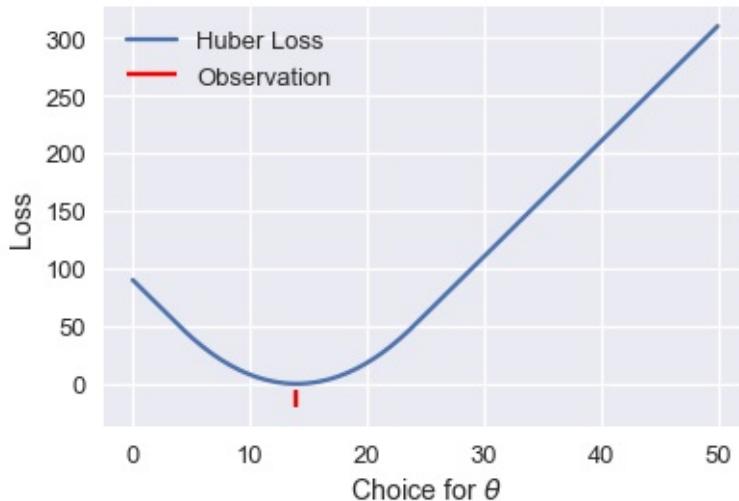


We can see that the Huber cost is smooth, unlike the mean absolute cost. The Huber cost also increases at a linear rate, unlike the quadratic rate of the mean squared cost.

The Huber cost does have a drawback, however. Notice that it transitions from the MSE to the mean absolute cost once θ gets far enough from the point. We can tweak this "far enough" to get different cost curves. For example, we can make it transition once θ is just one unit away from the observation:



Or we can make it transition when θ is ten units away from the observation:



This choice results in a different cost curve and can thus result in different values of θ . If we want to use the Huber cost function, we have the additional task of setting this transition point to a suitable value.

The Huber cost function is defined mathematically as follows:

$$\text{L}_{\alpha}(\theta, y) = \frac{1}{n} \sum_{i=1}^n \begin{cases} \frac{1}{2}(y_i - \theta)^2 & |y_i - \theta| \leq \alpha \\ \alpha(|y_i - \theta| - \frac{1}{2}\alpha) & \text{otherwise} \end{cases}$$

It is more complex than the previous cost functions because it combines both MSE and mean absolute cost. The additional parameter α sets the point where the Huber loss transitions from the MSE cost to the absolute cost.

Attempting to take the derivative of the Huber loss function is tedious and does not result in an elegant result like the MSE and mean absolute loss. Instead, we can use a computational method called gradient descent to find minimizing value of θ .

Gradient Descent and Numerical Optimization

In order to use a dataset for estimation and prediction, we need to precisely define our model and select a cost function. For example, in the tip percentage dataset, our model assumed that there was a single tip percentage that does not vary by table. Then, we decided to use the mean squared error cost function and found the model that minimized the cost function.

We also found that there are simple expressions that minimize the MSE and the mean absolute cost functions: the mean and the median. However, as our models and cost functions become more complicated we will no longer be able to find useful algebraic expressions for the models that minimize the cost. For example, the Huber cost has useful properties but is difficult to differentiate by hand.

We can use the computer to address this issue using gradient descent, a computational method of minimizing cost functions.

Show Widgets [Open on DataHub](#)

Table of Contents

- [Cost Minimization Using a Program](#)
- [Issues with `simple_minimize`](#)

Cost Minimization Using a Program

Let us return to our constant model:

$$\theta = C$$

We will use the mean squared error cost function:

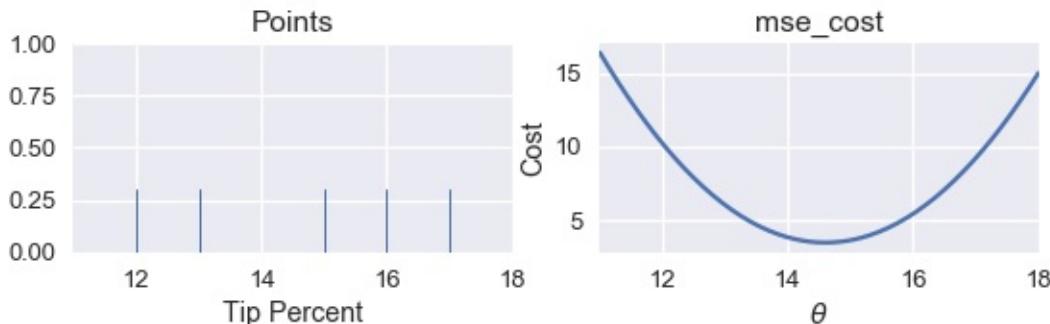
$\text{L}(\theta, y) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2$

For simplicity, we will use the dataset $y = [12, 13, 15, 16, 17]$. We know from our analytical approach in a previous chapter that the minimizing θ for the MSE cost is $\text{mean}(y) = 14.6$. Let's see whether we can find the same value by writing a program.

If we write the program well, we will be able to use the same program on any cost function in order to find the minimizing value of θ , including the mathematically complicated Huber cost:

$$\text{L}_\alpha(\theta, y) = \frac{1}{n} \sum_{i=1}^n \begin{cases} \frac{1}{2} (y_i - \theta)^2 & |y_i - \theta| \leq \alpha \\ \alpha (|y_i - \theta| - \frac{1}{2} \alpha) & \text{otherwise} \end{cases}$$

First, we create a rug plot of the data points. To the right of the rug plot we plot the MSE cost for different values of θ .



How might we write a program to automatically find the minimizing value of θ ? The simplest method is to compute the cost for many values θ . Then, we can return the θ value that resulted in the least cost.

We define a function called `simple_minimize` that takes in a cost function, an array of data points, and an array of θ values to try.

```
def simple_minimize(cost_fn, dataset, thetas):
    """
    Returns the value of theta in thetas that produces the least
    cost
    on a given dataset.
    """
    costs = [cost_fn(theta, dataset) for theta in thetas]
    return thetas[np.argmin(costs)]
```

Then, we can define a function to compute the MSE cost and pass it into `simple_minimize`.

```
def mse_cost(theta, dataset):
    return np.mean((dataset - theta) ** 2)

dataset = np.array([12, 13, 15, 16, 17])
thetas = np.arange(12, 18, 0.1)

simple_minimize(mse_cost, dataset, thetas)
```

14.599999999999991

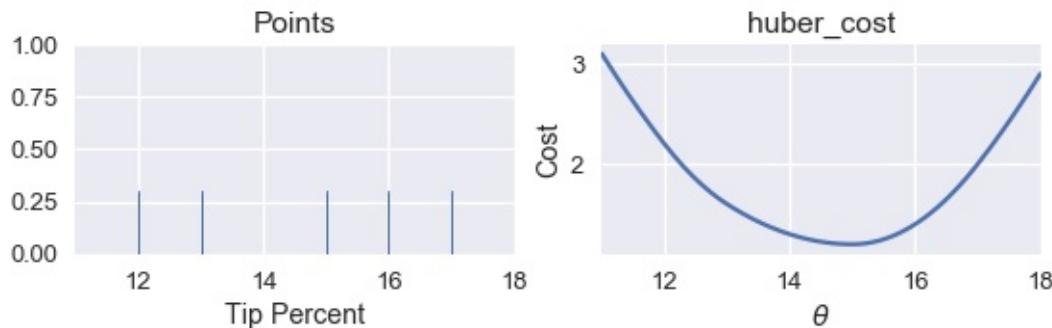
This is close to the expected value:

```
# Compute the minimizing theta using the analytical formula
np.mean(dataset)
```

14.6

Now, we can define a function to compute the Huber cost and plot the cost against θ .

```
def huber_cost(theta, dataset, alpha = 1):
    d = np.abs(theta - dataset)
    return np.mean(
        np.where(d < alpha,
                 (theta - dataset)**2 / 2.0,
                 alpha * (d - alpha / 2.0)))
)
```



Although we can see that the minimizing value of θ should be close to 15, we do not have an analytical method of finding θ directly for the Huber cost. Instead, we can use our `simple_minimize` function.

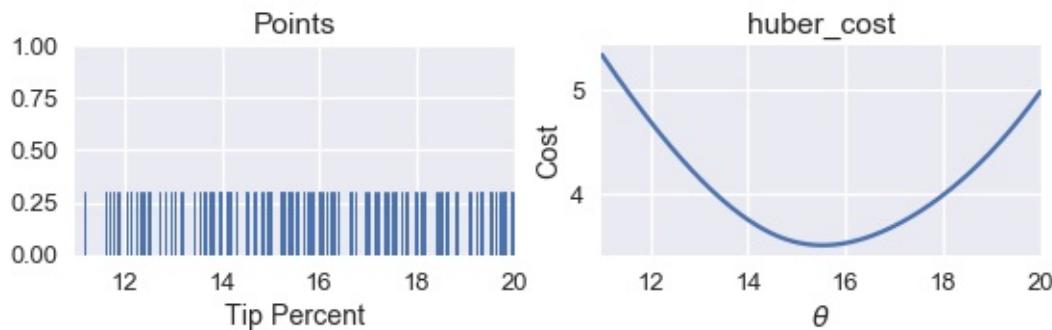
```
simple_minimize(huber_cost, dataset, thetas)
```

```
14.99999999999989
```

Now, we can return to our original dataset of tip percentages and find the best value for θ using the Huber cost.

```
tips = sns.load_dataset('tips')
tips['pcttip'] = tips['tip'] / tips['total_bill'] * 100
tips.head()
```

	total_bill	tip	sex	smoker	day	time	size	pcttip
0	16.99	1.01	Female	No	Sun	Dinner	2	5.944673
1	10.34	1.66	Male	No	Sun	Dinner	3	16.054159
2	21.01	3.50	Male	No	Sun	Dinner	3	16.658734
3	23.68	3.31	Male	No	Sun	Dinner	2	13.978041
4	24.59	3.61	Female	No	Sun	Dinner	4	14.680765



```
simple_minimize(huber_cost, tips['pcttip'], thetas)
```

```
15.49999999999988
```

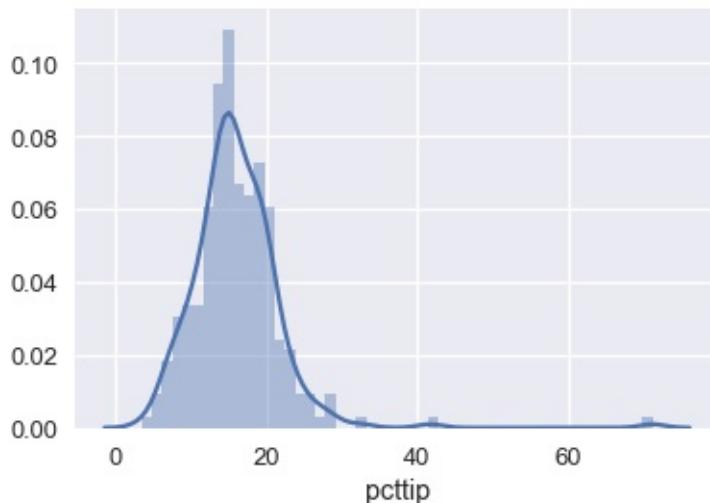
We can see that using the Huber cost gives us $\theta = 15.5$. We can now compare the minimizing θ values for MSE cost, mean absolute cost, and Huber cost.

```
print(f"      MSE cost: theta =\n{tips['pcttip'].mean():.2f}")
print(f"Mean Absolute cost: theta =\n{tips['pcttip'].median():.2f}")
print(f"      Huber cost: theta = 15.50")
```

```
MSE cost: theta = 16.08
Mean Absolute cost: theta = 15.48
Huber cost: theta = 15.50
```

We can see that the Huber cost is closer to the mean absolute cost since it is less affected by the outliers on the right side of the tip percentage distribution:

```
sns.distplot(tips['pcttip'], bins=50);
```



Issues with `simple_minimize`

Although `simple_minimize` allows us to minimize cost functions, it has some flaws that make it unsuitable for general purpose use. Its primary issue is that it only works with predetermined values of θ to test. For example, in this code snippet we used above, we had to manually define θ values in between 12 and 18.

```
dataset = np.array([12, 13, 15, 16, 17])
thetas = np.arange(12, 18, 0.1)

simple_minimize(mse_cost, dataset, thetas)
```

How did we know to examine the range between 12 and 18? We had to inspect the plot of the cost function manually and see that there was a minima in that range. This process becomes impractical as we add extra complexity to our models. In addition, we manually specified a step size of 0.1 in the code above. However, if the optimal value of θ was 12.043, our `simple_minimize` function would round to 12.00, the nearest multiple of 0.1.

We can solve both of these issues at once by using a method called *gradient descent*.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Gradient Descent](#)
 - [Intuition](#)
 - [Gradient Descent Analysis](#)
 - [Defining the `minimize` Function](#)
 - [Minimizing the Huber cost](#)
 - [Summary](#)

Gradient Descent

We are interested in creating a function that can minimize a cost function without forcing the user to predetermine which values of θ to try. In other words, while the `simple_minimize` function has the following signature:

```
simple_minimize(cost_fn, dataset, thetas)
```

We would like a function that has the following signature:

```
minimize(cost_fn, dataset)
```

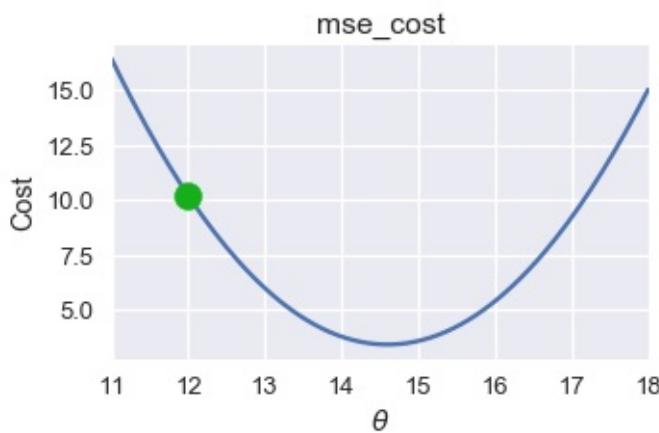
This function needs to automatically find the minimizing θ value no matter how small or large it is. We will use a technique called gradient descent to implement this new `minimize` function.

Intuition

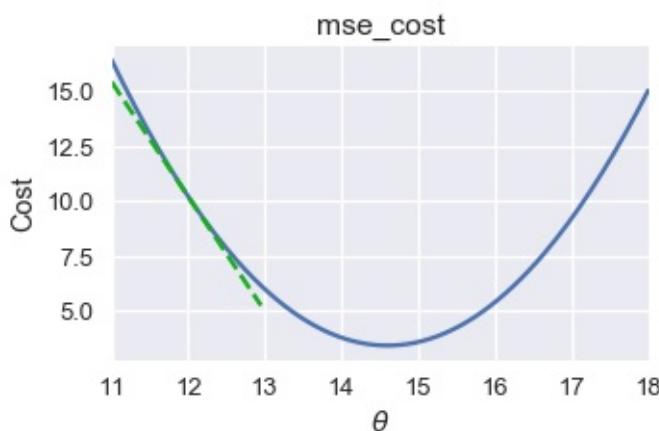
As with cost functions, we will discuss the intuition for gradient descent first, then formalize our understanding with mathematics.

Since the `minimize` function isn't given values of θ to try, we start by picking a θ anywhere we'd like. Then, we can iteratively improve the estimate of θ . To improve an estimate of θ , we look at the slope of the cost function at that choice of θ .

For example, suppose we are using MSE cost for the simple dataset $y = [12, 13, 15, 16, 17]$ and our current choice of θ is 12.

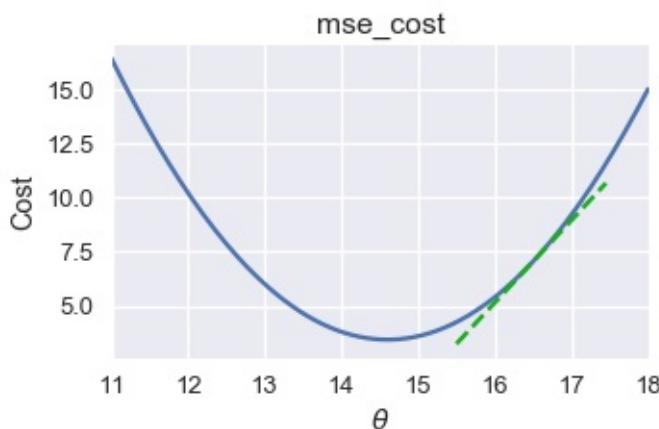


We'd like to choose a new value for θ that decreases the cost. To do this, we look at the slope of the cost function at $\theta = 12$:



The slope is negative, which means that increasing θ will decrease the cost.

If $\theta = 16.5$ on the other hand, the slope would be positive:



Which means that decreasing θ will decrease the cost.

Thus, the slope of the tangent line tells us which direction to move θ in order to decrease the cost. If the slope is negative, we want θ to move in the positive direction. If the slope is positive, θ should move in the negative direction.

Mathematically, we write:

$\$ \theta_{t+1} = \theta_t - \frac{\partial}{\partial \theta} L(\theta, y) \$\$$

Where θ_t is the current estimate and θ_{t+1} is the next estimate.

For the MSE cost, we have:

$$\$ \begin{aligned} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n (y_i - \theta)^2 \\ \frac{\partial}{\partial \theta} L(\theta, y) &= \frac{1}{n} \sum_{i=1}^n -2(y_i - \theta) \\ &= -\frac{2}{n} \sum_{i=1}^n (y_i - \theta) \end{aligned} \$\$$$

When $\theta_t = 12$, we can compute $-\frac{2}{n} \sum_{i=1}^n (y_i - \theta) = -5.2$.

Thus, $\theta_{t+1} = 12 - (-5.2) = 17.2$.

We've plotted the old value of θ as a green outlined circle and the new value as a filled in circle on the cost curve below.



Although θ went in the right direction, it ended up as far away from the minimum as it started. We can remedy this by multiplying the slope by a small constant before subtracting it from θ . Our final update formula is:

$\$ \theta_{t+1} = \theta_t - \alpha \cdot \frac{\partial}{\partial \theta} L(\theta, y) \$\$$

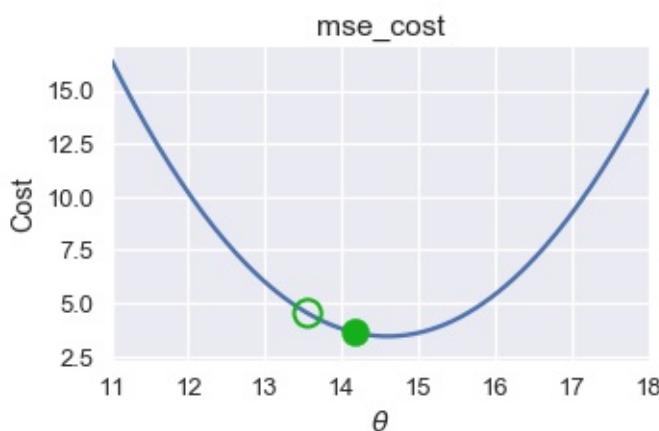
Where α is a small constant. For example, if we set $\alpha = 0.3$ this is the new θ_{t+1} :

```
old theta: 12
new theta: 13.56
```

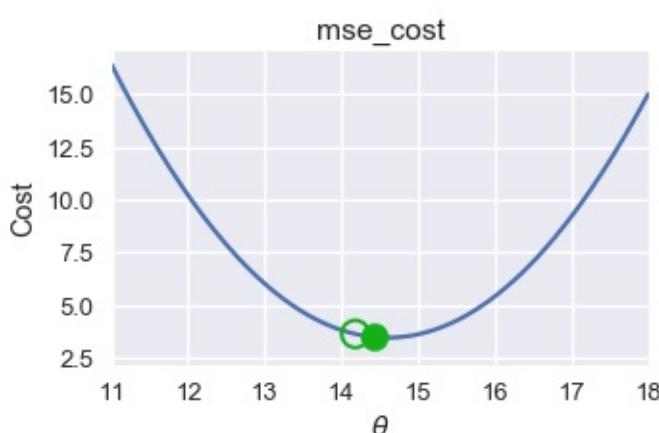


Here are the θ values for successive iterations of this process. Notice that θ changes more slowly as it gets closer to the minimum cost because the slope is also smaller.

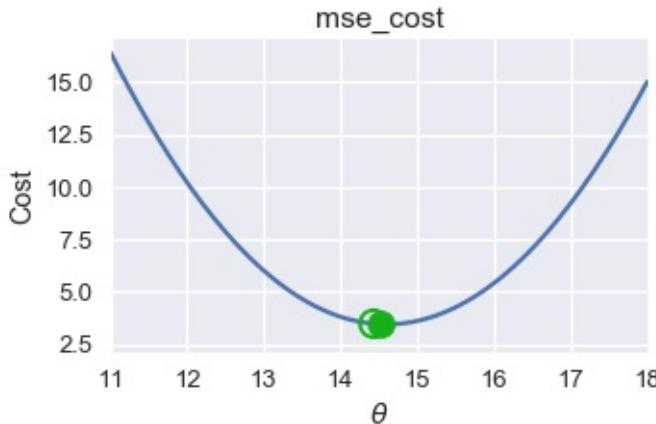
```
old theta: 13.56
new theta: 14.184000000000001
```



```
old theta: 14.18
new theta: 14.432
```



```
old theta: 14.432
new theta: 14.5328
```



Gradient Descent Analysis

We now have the full algorithm for gradient descent:

1. Choose a starting value of θ (0 is a common choice).
2. Compute $\theta - \alpha \frac{\partial}{\partial \theta} L(\theta, y)$ and store this as the new value of θ .
3. Repeat until θ doesn't change between iterations.

You will more commonly see the gradient ∇_{θ} in place of the partial derivative $\frac{\partial}{\partial \theta}$. The two notations are essentially equivalent, but since the gradient notation is more common we will use it in the gradient update formula from now on:

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} L(\theta, y) \quad \text{--- (1)}$$

To review notation:

- θ_t is the current choice of θ .
- θ_{t+1} is the next choice of θ .
- α is called the learning rate, usually set to a small constant. Sometimes it is useful to start with a larger α and decrease it over time. If α changes between iterations, we use the variable α_t to mark that α varies over time t .
- $\nabla_{\theta} L(\theta, y)$ is the partial derivative / gradient of the cost function at θ .

You can now see the importance of choosing a differentiable cost function: $\nabla_{\theta} L(\theta, y)$ is a crucial part of the gradient descent algorithm. (While it is possible to estimate the gradient by computing the difference in cost for two slightly different values of θ ,

\theta \$ and dividing by the distance between \$ \theta \$ values, this typically increases the runtime of gradient descent so significantly that it becomes impractical to use.)

The gradient algorithm is simple yet powerful since we can use it for many types of models and many types of cost functions. It is the computational tool of choice for fitting many important models, including linear regression on large datasets and neural networks.

Defining the `minimize` Function¶

Now we return to our original task: defining the `minimize` function. We will have to change our function signature slightly since we now need to compute the gradient of the cost function.

```
def minimize(cost_fn, grad_cost_fn, dataset, alpha=0.2,
progress=True):
    """
    Uses gradient descent to minimize cost_fn. Returns the
    minimizing value of
    theta once theta changes less than 0.001 between iterations.
    """
    theta = 0
    while True:
        if progress:
            print(f'theta: {theta:.2f} | cost: {cost_fn(theta,
dataset):.2f}')
        gradient = grad_cost_fn(theta, dataset)
        new_theta = theta - alpha * gradient

        if abs(new_theta - theta) < 0.001:
            return new_theta

        theta = new_theta
```

Then we can define functions to compute our MSE cost and its gradient:

```
def mse_cost(theta, y_vals):
    return np.mean((y_vals - theta) ** 2)

def grad_mse_cost(theta, y_vals):
    return -2 * np.mean(y_vals - theta)
```

Finally, we can use the `minimize` function to compute the minimizing value of θ for $y = [12, 13, 15, 16, 17]$.

```
%%time
theta = minimize(mse_cost, grad_mse_cost, np.array([12, 13, 15,
16, 17]))
print(f'Minimizing theta: {theta}')
print()
```

```
theta: 0.00 | cost: 216.60
theta: 5.84 | cost: 80.18
theta: 9.34 | cost: 31.07
theta: 11.45 | cost: 13.39
theta: 12.71 | cost: 7.02
theta: 13.46 | cost: 4.73
theta: 13.92 | cost: 3.90
theta: 14.19 | cost: 3.61
theta: 14.35 | cost: 3.50
theta: 14.45 | cost: 3.46
theta: 14.51 | cost: 3.45
theta: 14.55 | cost: 3.44
theta: 14.57 | cost: 3.44
theta: 14.58 | cost: 3.44
theta: 14.59 | cost: 3.44
theta: 14.60 | cost: 3.44
theta: 14.60 | cost: 3.44
Minimizing theta: 14.59851722463264

CPU times: user 2.6 ms, sys: 1.6 ms, total: 4.2 ms
Wall time: 3.02 ms
```

We can see that gradient quickly finds the same solution as the analytic method:

```
np.mean([12, 13, 15, 16, 17])
```

```
14.6
```

Minimizing the Huber cost

Now, we can apply gradient descent to minimize the Huber cost on our dataset of tip percentages.

The Huber cost is:

$$\$ L_{\delta}(\theta, y) = \frac{1}{n} \sum_{i=1}^n \begin{cases} \frac{1}{2}(y_i - \theta)^2 & |y_i - \theta| \leq \delta \\ \delta(|y_i - \theta| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

The gradient of the Huber cost is:

$$\$ \nabla_{\theta} L_{\delta}(\theta, y) = \frac{1}{n} \sum_{i=1}^n \begin{cases} -(y_i - \theta) & |y_i - \theta| \leq \delta \\ -\delta \cdot \text{sign}(y_i - \theta) & \text{otherwise} \end{cases}$$

(Note that in previous definitions of Huber cost we used the variable α to denote the transition point. To avoid confusion with the α used in gradient descent, we replace the transition point parameter of the Huber loss with δ .)

```
def huber_cost(theta, dataset, delta = 1):
    d = np.abs(theta - dataset)
    return np.mean(
        np.where(d <= delta,
                  (theta - dataset)**2 / 2.0,
                  delta * (d - delta / 2.0)))
)

def grad_huber_cost(theta, dataset, delta = 1):
    d = np.abs(theta - dataset)
    return np.mean(
        np.where(d <= delta,
                  -(dataset - theta),
                  -delta * np.sign(dataset - theta)))
)
```

Let's minimize the Huber cost on the tips dataset:

```
%%time
theta = minimize(huber_cost, grad_huber_cost, tips['pcttip'],
progress=False)
print(f'Minimizing theta: {theta}')
print()
```

```
Minimizing theta: 15.506849531471964
```

```
CPU times: user 207 ms, sys: 6.9 ms, total: 213 ms
Wall time: 258 ms
```

Summary

Gradient descent gives us a generic way to minimize a cost function when we cannot solve for the minimizing value of θ analytically. As our models and cost functions increase in complexity, we will turn to gradient descent again and again as our tool of choice for model fitting.

Linear Models

Now that we have a general method for fitting a model to a cost function, we turn our attention to improvements on our model. For the sake of simplicity, we previously restricted ourselves to a constant model: our model only ever predicts a single number.

However, giving our waiter such a model would hardly satisfy him. He would likely point out that he collected much more information about his tables than simply the tip percents. Why didn't we use his other data—e.g. size of the table or total bill—in order to make our model more useful?

In this chapter we will introduce linear models which will allow us to make use of our entire dataset to make predictions. Linear models are not only widely used in practice but also have rich theoretical underpinnings that will allow us to understand future tools for modeling.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Predicting Tip Amounts](#)
- [Defining a Linear Model](#)
 - [Estimating the Linear Model](#)

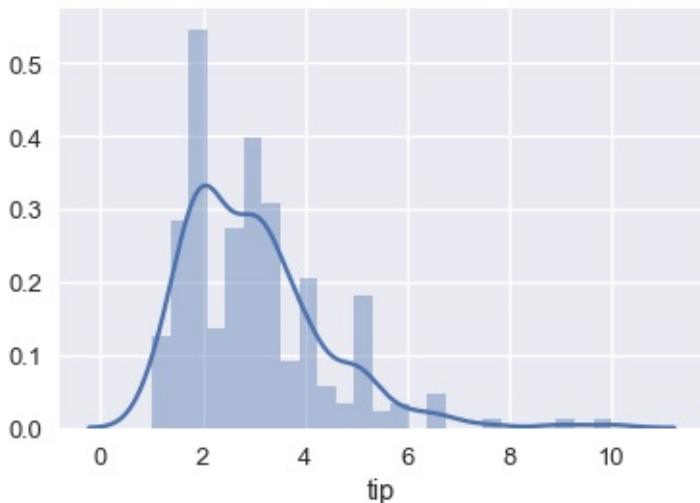
Predicting Tip Amounts

Previously, we worked with a dataset that contained one row for each table that a waiter served in a week. Our waiter collected this data in order to predict the tip amount he could expect to receive from a future table.

```
tips = sns.load_dataset('tips')
tips.head()
```

	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

```
sns.distplot(tips['tip'], bins=25);
```



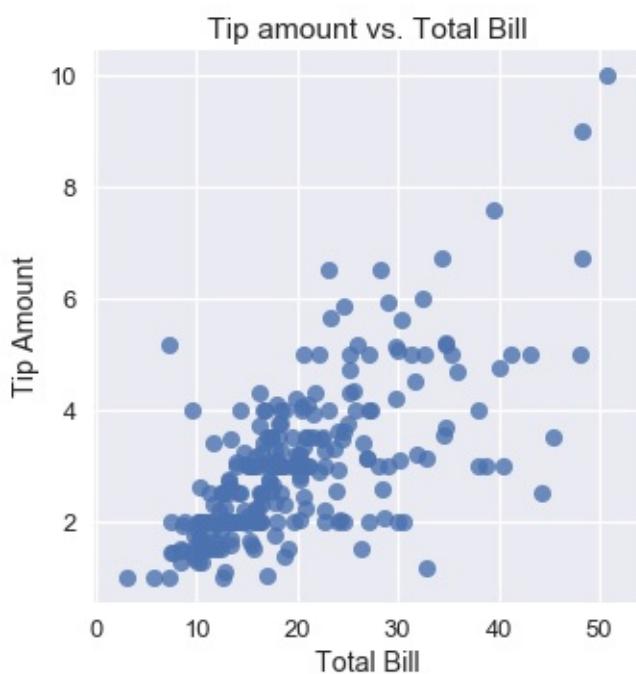
As we have covered previously, if we choose a constant model and the mean squared error cost our model will predict the mean of the tip amount:

```
np.mean(tips['tip'])
```

```
2.9982786885245902
```

This means that if a new party orders a meal and the waiter asks us how much tip he will likely receive, we will say "around \\$3", no matter how large the table is or how much their total bill was.

However, looking at other variables in the dataset suggest that we might be able to make more accurate predictions if we incorporate them into our model. For example, the following plot of the tip amount against the total bill shows a positive association.



Although the average tip amount is \$3, if a table orders \$40 worth of food we would certainly expect that the waiter receives more than \$3 of tip. Thus, we would like to alter our model so that it makes predictions based on the variables in our dataset instead of blindly predicting the mean tip amount. To do this, we use a linear model instead of constant one.

Let's briefly review our current toolbox for modeling and estimation and define some new notation so that we can better represent the additional complexity that linear models have.

Defining a Linear Model

We are interested in predicting the tip amount based on the total bill of a table. We use y to represent the tip amount, the variable we are trying to predict. We use x to represent the total bill, the variable we use for prediction. Our model assumes that in truth, y has a perfectly linear relationship with x . However, our observed data do not follow a perfectly straight line because of some random noise ϵ . Mathematically, we write:

$$y = f_{\theta}(x) + \epsilon$$

where

$$f_{\theta}(x) = \theta_1 x + \theta_0$$

This means that if the assumptions of our model hold, we are able to somehow find the exact values of θ_1 and θ_0 , and we magically have no random noise, we will be able to perfectly predict the amount of tip the waiter will get for all tables, forever. Of course, we cannot completely fulfill any of these criteria in practice. Instead, we will estimate θ_1 and θ_0 using our dataset.

(Note that we use $f_{\theta}(x)$ to signify the true mechanism by which our data are generated. For example, our constant model assumed that the population tip amounts did not depend on x so for the constant model we would write $f_{\theta}(x) = \theta$.)

Estimating the Linear Model

Since we cannot find θ_1 and θ_0 exactly, we will assume that our dataset approximates our population and use our dataset to estimate these parameters. We denote our estimations with $\hat{\theta}_1$ and $\hat{\theta}_0$ and define our model as:

$$\hat{f}_{\theta}(x) = \hat{\theta}_1 x + \hat{\theta}_0$$

Sometimes you will see $h(x)$ written instead of $\hat{f}_{\theta}(x)$; the " h " stands for hypothesis.

In order to determine $\hat{\theta}_1$ and $\hat{\theta}_0$, we choose a cost function and minimize it using gradient descent.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Fitting a Linear Model Using Gradient Descent](#)
- [Derivative of the MSE cost](#)
- [Running Gradient Descent](#)

Fitting a Linear Model Using Gradient Descent

We want to fit a linear model that predicts the tip amount based on the total bill of the table:

$$\hat{f}(\theta)(x) = \hat{\theta}_1 x + \hat{\theta}_0$$

In order to find $\hat{\theta}_1$ and $\hat{\theta}_0$, we need to first choose a cost function. We will choose the mean squared error cost function:

$$L(\hat{\theta}, x, y) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(\theta)(x_i))^2$$

Note that we have modified our loss function to incorporate our new model. In this case x is a vector containing the individual total bills, y is a vector containing the individual tip amounts, and $\hat{\theta}$ is a vector: $\theta = [\hat{\theta}_1, \hat{\theta}_0]$.

Using a linear model with the squared error also goes by the name of least-squares linear regression. We can use gradient descent to find the θ that minimizes the cost.

An Aside on Using Correlation

If you have seen least-squares linear regression before, you may recognize that we can compute the correlation coefficient and use it to determine $\hat{\theta}_1$ and $\hat{\theta}_0$. This is simpler and faster to compute than using gradient descent for this specific problem, similar to how computing the mean was simpler than using gradient descent to fit a constant model. We will use gradient descent anyway because it is a general-purpose method of cost minimization that still works when we later introduce models that do not have analytic solutions. In fact, in many real-world scenarios we will use gradient descent even when we have an analytic solution because the analytic solution is too difficult to compute.

Derivative of the MSE cost

In order to use gradient descent, we have to compute the derivative of the MSE cost with respect to $\hat{\theta}$. Now that $\hat{\theta}$ is a vector of length 2 instead of a scalar, $\nabla_{\hat{\theta}} L(\hat{\theta}, x, y)$ will also be a vector of length 2.

$$\begin{aligned} \nabla_{\hat{\theta}} L(\hat{\theta}, x, y) &= \nabla_{\hat{\theta}} \left[\frac{1}{n} \sum_{i=1}^n (y_i - f_{\hat{\theta}}(x_i))^2 \right] \\ &= \frac{1}{n} \sum_{i=1}^n 2(y_i - f_{\hat{\theta}}(x_i))(-\nabla_{\hat{\theta}} f_{\hat{\theta}}(x_i)) \\ &= -\frac{2}{n} \sum_{i=1}^n (y_i - f_{\hat{\theta}}(x_i)) \nabla_{\hat{\theta}} f_{\hat{\theta}}(x_i) \end{aligned}$$

We know:

$$f_{\hat{\theta}}(x) = \hat{\theta}_0 + \hat{\theta}_1 x$$

We now need to compute $\nabla_{\hat{\theta}} f_{\hat{\theta}}(x_i)$ which is a length 2 vector.

$$\begin{aligned} \nabla_{\hat{\theta}} f_{\hat{\theta}}(x_i) &= \begin{pmatrix} \frac{\partial}{\partial \hat{\theta}_0} f_{\hat{\theta}}(x_i) & \frac{\partial}{\partial \hat{\theta}_1} f_{\hat{\theta}}(x_i) \end{pmatrix} \\ &= \begin{pmatrix} \frac{\partial}{\partial \hat{\theta}_0} (\hat{\theta}_0 + \hat{\theta}_1 x_i) & \frac{\partial}{\partial \hat{\theta}_1} (\hat{\theta}_0 + \hat{\theta}_1 x_i) \end{pmatrix} \\ &= \begin{pmatrix} 1 & x_i \end{pmatrix} \end{aligned}$$

Finally, we plug back into our formula above to get

$$\begin{aligned} \nabla_{\hat{\theta}} L(\hat{\theta}, x, y) &= -\frac{2}{n} \sum_{i=1}^n (y_i - f_{\hat{\theta}}(x_i)) \begin{pmatrix} 1 & x_i \end{pmatrix} \\ &= -\frac{2}{n} \sum_{i=1}^n (y_i - f_{\hat{\theta}}(x_i)) \begin{pmatrix} 1 & x_i \end{pmatrix} (y_i - f_{\hat{\theta}}(x_i)) \end{aligned}$$

This is a length 2 vector since $(y_i - f_{\hat{\theta}}(x_i))$ is a scalar.

Running Gradient Descent

Now, let's fit a linear model on the tips dataset in to predict the tip amount from the total table bill. First, we define a Python function to compute the cost:

```
def simple_linear_model(thetas, x_vals):
    '''Returns predictions by a linear model on x_vals.'''
    return thetas[0] + thetas[1] * x_vals

def mse_cost(thetas, x_vals, y_vals):
    return np.mean((y_vals - simple_linear_model(thetas,
        x_vals)) ** 2)
```

Then, we define a function to compute the gradient of the cost:

```
def grad_mse_cost(thetas, x_vals, y_vals):
    n = len(x_vals)
    grad_0 = y_vals - simple_linear_model(thetas, x_vals)
    grad_1 = (y_vals - simple_linear_model(thetas, x_vals)) *
    x_vals
    return -2 / n * np.array([np.sum(grad_0), np.sum(grad_1)])
```

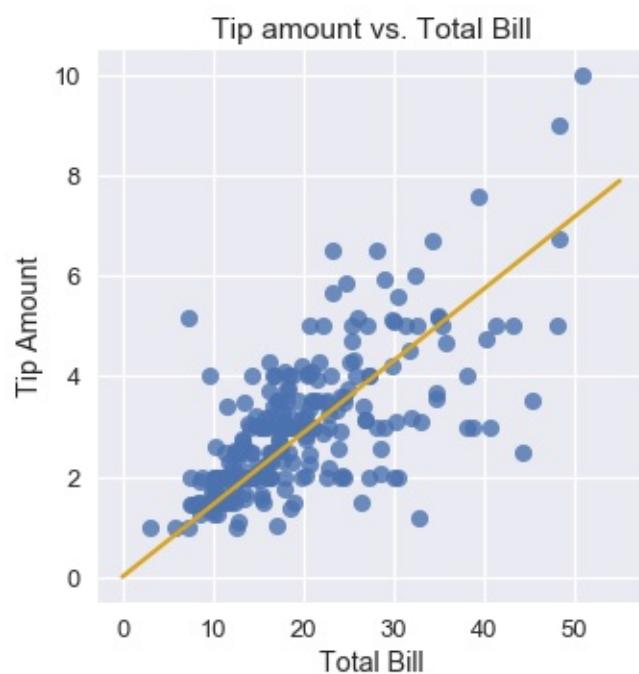
Finally, we run gradient descent!

```
%%time

thetas = minimize(mse_cost, grad_mse_cost, tips['total_bill'],
tips['tip'])
```

```
theta: [ 0.  0.] | cost: 10.896283606557377
theta: [ 0.      0.07] | cost: 3.8937622006094705
theta: [ 0.      0.1] | cost: 1.9359443267168215
theta: [ 0.01   0.12] | cost: 1.388538448286097
theta: [ 0.01   0.13] | cost: 1.235459416905535
theta: [ 0.01   0.14] | cost: 1.1926273731479433
theta: [ 0.01   0.14] | cost: 1.1806184944517062
theta: [ 0.01   0.14] | cost: 1.177227251696266
theta: [ 0.01   0.14] | cost: 1.1762453624313751
theta: [ 0.01   0.14] | cost: 1.1759370980989148
theta: [ 0.01   0.14] | cost: 1.175817178966766
CPU times: user 57.7 ms, sys: 14.1 ms, total: 71.8 ms
Wall time: 76.9 ms
```

We can plot our predictions on the scatter plot of tip vs. total bill.



We can see that if a table's bill is \\$10, our model will predict that the waiter gets around \\$1.50 in tip. Similarly, if a table's bill is \\$40, our model will predict a tip of around \\$6.00.

[Show Widgets](#) [Open on DataHub](#)

Table of Contents

- [Multiple Linear Regression](#)
- [MSE Cost and its Gradient](#)
- [Fitting the Model With Gradient Descent](#)
- [Visualizing our Predictions](#)
- [Using All the Data](#)
- [Summary](#)

Multiple Linear Regression

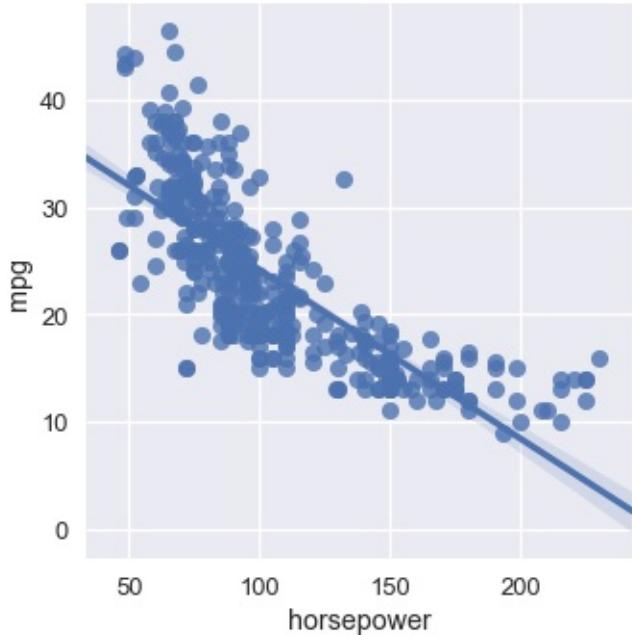
Our simple linear model has a key advantage over the constant model: it uses the data when making predictions. However, it is still rather limited since it can only use one variable in our dataset when many datasets have many potentially useful variables. For example, consider the following dataset on car models and their milage per gallon (MPG):

```
mpg = pd.read_csv('mpg.csv').dropna().reset_index(drop=True)
mpg
```

	mpg	cylinders	displacement	...	model year	origin	car name
0	18.0	8	307.0	...	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	...	70	1	buick skylark 320
2	18.0	8	318.0	...	70	1	plymouth satellite
...
389	32.0	4	135.0	...	82	1	dodge rampage
390	28.0	4	120.0	...	82	1	ford ranger
391	31.0	4	119.0	...	82	1	chevy s-10

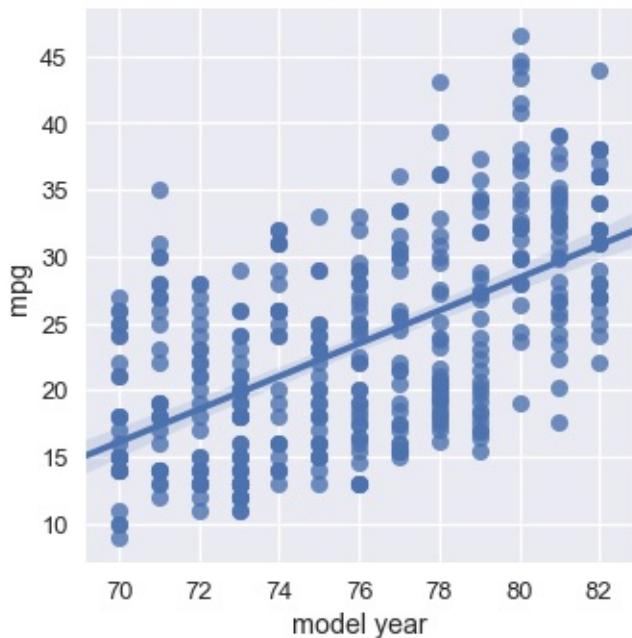
392 rows × 9 columns

It seems likely that multiple attributes of a car model will affect its MPG. For example, the MPG seems to decrease as horsepower increases:



However, cars released later generally have better MPG than older cars:

```
sns.lmplot(x='model year', y='mpg', data=mpg);
```



It seems possible that we can get a more accurate model if we could take both horsepower and model year into account when making predictions about the MPG. In fact, perhaps the best model takes into account all the numerical variables in our dataset. We can extend our univariate linear regression to allow prediction based on any number of attributes.

We state the following model:

$\hat{f}(\theta)(x) = \hat{\theta}_0 + \hat{\theta}_1 x_1 + \dots + \hat{\theta}_p x_p$

Where x now represents a vector containing p attributes of a single car. The model above says, "Take multiple attributes of a car, multiply them by some weights, and add them together to make a prediction for MPG."

For example, if we're making a prediction on the first car in our dataset using the horsepower, weight, and model year columns, the vector x looks like:

	horsepower	weight	model year
0	130.0	3504.0	70

In these examples, we've kept the column names for clarity but keep in mind that x only contains the numerical values of the table above: $x = [130.0, 3504.0, 70]$.

Now, we will perform a notational trick that will greatly simplify later formulas. We will prepend the value 1 to the vector x , so that we have the following vector for x :

	bias	horsepower	weight	model year
0	1	130.0	3504.0	70

Now, observe what happens to the formula for our model:

$$\begin{aligned} \hat{f}(\theta)(x) &= \hat{\theta}_0 + \hat{\theta}_1 x_1 + \dots + \hat{\theta}_p x_p \\ &= \hat{\theta}_0 (1) + \hat{\theta}_1 x_1 + \dots + \hat{\theta}_p x_p \\ &= \hat{\theta}_0 x_0 + \hat{\theta}_1 x_1 + \dots + \hat{\theta}_p x_p \end{aligned}$$

Where $\hat{\theta} \cdot x$ is the vector dot product of $\hat{\theta}$ and x . Vector and matrix notation was designed to succinctly write linear combinations and is thus well-suited for our linear models. However, you will have to remember from now on that $\hat{\theta} \cdot x$ is a vector-vector dot product. When in doubt, you can always expand the dot product into simple multiplications and additions.

Now, we define the matrix X as the matrix containing every car model as a row. For example, here are the first five rows of X :

	bias	horsepower	weight	model year
0	1	130.0	3504.0	70
1	1	165.0	3693.0	70
2	1	150.0	3436.0	70
3	1	150.0	3433.0	70
4	1	140.0	3449.0	70

Again, keep in mind that the actual matrix X only contains the numerical values of the table above.

Notice that X is composed of multiple x vectors stacked on top of each other. To keep the notation clear, we define X_i to refer to the row with index i of X . We define X_{ij} to refer to the element with index j of the row with index i of X . Thus, X_i is a p -dimensional vector and X_{ij} is a scalar. X is an $n \times p$ matrix.

For example, from the table above we have $X_4 = [1, 140, 3449, 70]$ and $X_{41} = 140$. This notation becomes important when we define the cost function since we will need both X , the matrix of input values, and y , the vector of MPGs.

MSE Cost and its Gradient

The mean squared error cost function takes in a vector of weights $\hat{\theta}$, a matrix of inputs X , and a vector of observed MPGs y :

$$\begin{aligned} L(\hat{\theta}, X, y) &= \frac{1}{n} \sum_i (y_i - f(\hat{\theta})(X_i))^2 \end{aligned}$$

We've previously derived the gradient of the MSE cost with respect to $\hat{\theta}$:

$$\begin{aligned} \nabla_{\hat{\theta}} L(\hat{\theta}, X, y) &= -\frac{2}{n} \sum_i (y_i - f(\hat{\theta})(X_i)) \nabla_{\hat{\theta}} f(\hat{\theta})(X_i) \end{aligned}$$

We know that:

$$f(\hat{\theta})(x) = \hat{\theta}_0 x_0 + \dots + \hat{\theta}_p x_p$$

Let's now compute $\nabla_{\hat{\theta}} f(\hat{\theta})(x)$. The result is surprisingly simple because $\hat{\theta} \cdot x = \hat{\theta}_0 x_0 + \dots + \hat{\theta}_p x_p$ and thus $\frac{\partial}{\partial \hat{\theta}_0} (\hat{\theta} \cdot x) = x_0$, $\frac{\partial}{\partial \hat{\theta}_1} (\hat{\theta} \cdot x) = x_1$, and so on.

$$\begin{aligned} \nabla_{\hat{\theta}} f(\hat{\theta})(x) &= \begin{bmatrix} \nabla_{\hat{\theta}_0} f(\hat{\theta})(x) & \nabla_{\hat{\theta}_1} f(\hat{\theta})(x) & \dots & \nabla_{\hat{\theta}_p} f(\hat{\theta})(x) \end{bmatrix} \\ &= \begin{bmatrix} x_0 & x_1 & \dots & x_p \end{bmatrix} \end{aligned}$$

Finally, we plug this result back into our gradient calculations:

$$\begin{aligned} \nabla_{\hat{\theta}} L(\hat{\theta}, X, y) &= -\frac{2}{n} \sum_i (y_i - f(\hat{\theta})(X_i)) \nabla_{\hat{\theta}} f(\hat{\theta})(X_i) \\ &= -\frac{2}{n} \sum_i (y_i - \hat{\theta} \cdot X_i) X_i \end{aligned}$$

Remember that since $y_i - \hat{\theta}^\top X_i$ is a scalar and X_i is a p -dimensional vector, the gradient $\nabla_{\hat{\theta}} L(\hat{\theta}, X, y)$ is a p -dimensional vector.

We saw this same type of result when we computed the gradient for univariate linear regression and found that it was 2-dimensional since θ was 2-dimensional.

Fitting the Model With Gradient Descent

We can now plug in our cost and its derivative into gradient descent. As usual, we will define the model, cost, and gradient cost in Python.

```
def linear_model(thetas, X):
    '''Returns predictions by a linear model on x_vals.'''
    return X @ thetas

def mse_cost(thetas, X, y):
    return np.mean((y - linear_model(thetas, X)) ** 2)

def grad_mse_cost(thetas, X, y):
    n = len(X)
    return -2 / n * (X.T @ y - X.T @ X @ thetas)
```

Now, we can simply plug in our functions into our gradient descent minimizer:

```
%%time

thetas = minimize(mse_cost, grad_mse_cost, X, y)
print(f'theta: {thetas} | cost: {mse_cost(thetas, X, y):.2f}')
```

```

theta: [ 0.  0.  0.  0.] | cost: 610.47
theta: [ 0.      0.      0.01  0.    ] | cost: 178.95
theta: [ 0.01 -0.11 -0.     0.55] | cost: 15.78
theta: [ 0.01 -0.01 -0.01  0.58] | cost: 11.97
theta: [-4.     -0.01 -0.01  0.63] | cost: 11.81
theta: [-13.72 -0.     -0.01  0.75] | cost: 11.65
theta: [-13.72 -0.     -0.01  0.75] | cost: 11.65
CPU times: user 8.81 ms, sys: 3.11 ms, total: 11.9 ms
Wall time: 9.22 ms

```

Visualizing our Predictions

How does our model do? We can see that the cost decreased dramatically (from 610 to 11.6). We can show the predictions of our model alongside the original values:

	predicted_mpg	mpg	horsepower	weight	model year
0	15.447125	18.0	130.0	3504.0	70
1	14.053509	15.0	165.0	3693.0	70
2	15.785576	18.0	150.0	3436.0	70
...
389	32.456900	32.0	84.0	2295.0	82
390	30.354143	28.0	79.0	2625.0	82
391	29.726608	31.0	82.0	2720.0	82

392 rows × 5 columns

Since we found $\hat{\theta}$ from gradient descent, we can verify for the first row of our data that $\hat{\theta} \cdot X_0$ matches our prediction above:

```

print(f'Prediction for first row: '
      f'{thetas[0] + thetas[1] * 130 + thetas[2] * 3504 +
thetas[3] * 70:.2f}')

```

```
Prediction for first row: 15.45
```

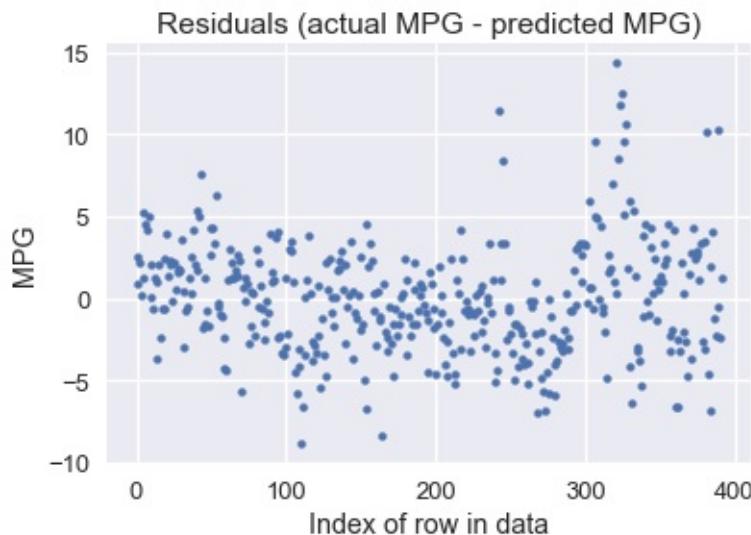
We've included a widget below to pan through the predictions and the data used to make the prediction:

Show Widget

(392 rows, 5 columns) total

We can also plot the residuals of our predictions (actual values - predicted values):

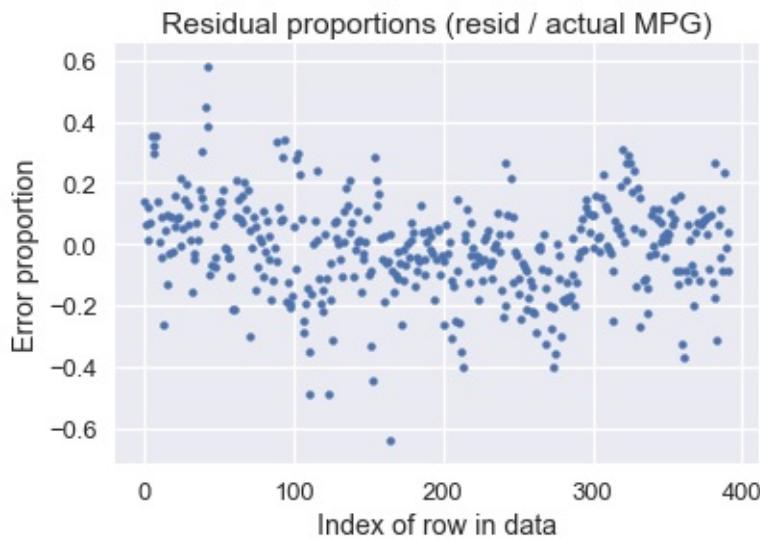
```
resid = y - linear_model(thetas, X)
plt.scatter(np.arange(len(resid)), resid, s=15)
plt.title('Residuals (actual MPG - predicted MPG)')
plt.xlabel('Index of row in data')
plt.ylabel('MPG');
```



It looks like our model makes reasonable predictions for many car models, although there are some predictions that were off by over 10 MPG (some cars had under 10 MPG!).

Perhaps we are more interested in the percent error between the predicted MPG values and the actual MPG values:

```
resid_prop = resid / with_predictions['mpg']
plt.scatter(np.arange(len(resid_prop)), resid_prop, s=15)
plt.title('Residual proportions (resid / actual MPG)')
plt.xlabel('Index of row in data')
plt.ylabel('Error proportion');
```



It looks like our model's predictions are usually around 20% away from the actual MPG values.

Using All the Data

Notice that in our example thus far, our X matrix has four columns: one column of all ones, the horsepower, the weight, and the model year. However, model allows us to handle an arbitrary number of columns:

$$\$ \$ \begin{aligned} f_{\hat{\theta}}(x) &= \hat{\theta}_0 + \hat{\theta}_1 x_1 + \hat{\theta}_2 x_2 + \dots + \hat{\theta}_n x_n \end{aligned} \$ \$$$

As we include more columns into our data matrix, we extend $\hat{\theta}$ so that it has one parameter for each column in X . Instead of only selecting three numerical columns for prediction, why not use all seven of them?

	bias	cylinders	displacement	...	acceleration	model year	origin
0	1	8	307.0	...	12.0	70	1
1	1	8	350.0	...	11.5	70	1
2	1	8	318.0	...	11.0	70	1
...
389	1	4	135.0	...	11.6	82	1
390	1	4	120.0	...	18.6	82	1
391	1	4	119.0	...	19.4	82	1

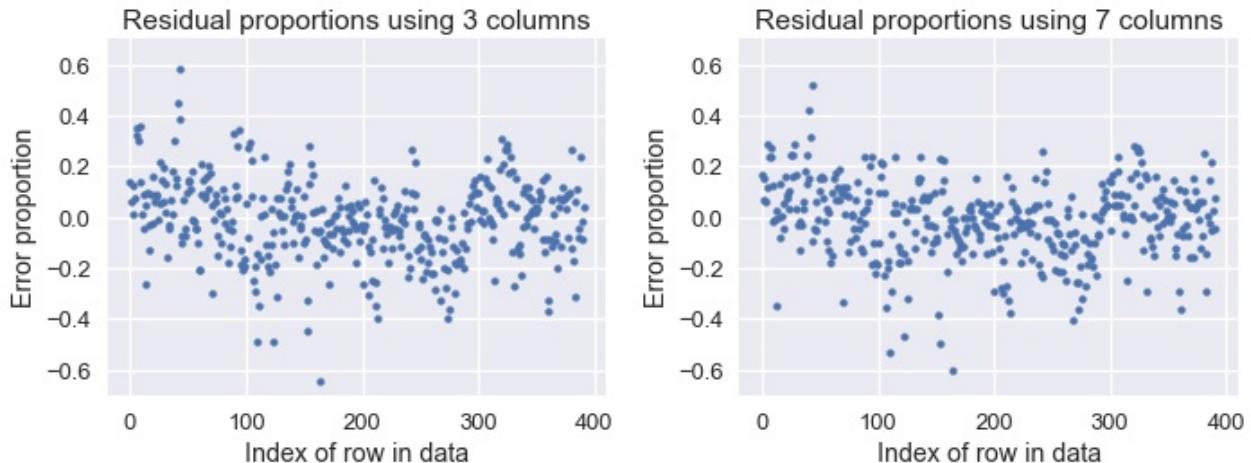
392 rows \times 8 columns

```
%%time

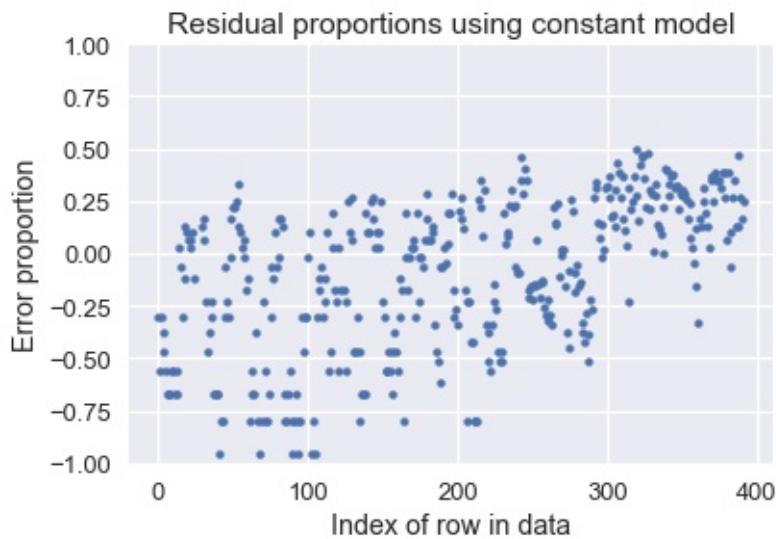
thetas_all = minimize(mse_cost, grad_mse_cost, X, y,
progress=10)
print(f'theta: {thetas_all} | cost: {mse_cost(thetas_all, X,
y):.2f}')
```

```
theta: [ 0.  0.  0.  0.  0.  0.  0.] | cost: 610.47
theta: [-0.5 -0.81  0.02 -0.04 -0.01 -0.07  0.59  1.3 ] | cost:
11.22
theta: [-17.23 -0.49   0.02 -0.02 -0.01   0.08   0.75   1.43]
| cost: 10.85
theta: [-17.22 -0.49   0.02 -0.02 -0.01   0.08   0.75   1.43]
| cost: 10.85
CPU times: user 10.9 ms, sys: 3.51 ms, total: 14.4 ms
Wall time: 11.7 ms
```

We see that our cost has decreased from 11.6 with three columns of our dataset to 10.85 when using all seven numerical columns of our dataset. We display the proportion error plots for both old and new predictions below:



Although the difference is slight, you can see that the errors are a bit lower when using seven columns compared to using three. Both models are much better than using a constant model, as the below plot shows:



Using a constant model results in over 75% error for many car MPGs!

Summary

We have introduced the linear model for regression. Unlike the constant model, the linear model takes in features of our data into account when making predictions, making it much more useful whenever we have correlations between variables of our data.

The procedure of fitting a model to data should now be quite familiar:

1. Select a model.
2. Select a cost function.
3. Minimize the cost function using gradient descent.

It is useful to know that we can usually tweak one of the components without changing the others. In this section, we introduced the linear model without changing our cost function or using a different minimization algorithm. Although modeling can get complicated, it is usually easier to learn by focusing on one component at a time, then combining different parts together as needed in practice.