

Module 3 - Statistics Essentials and Data Pre-processing with Python

Data Pre-processing with Pandas

In the previous notebook, we dove into detail on NumPy and its `ndarray` object, which provides efficient storage and manipulation of dense typed arrays in Python. Here we'll build on this knowledge by looking in detail at the data structures provided by the Pandas library. Pandas is a newer package built on top of NumPy, and provides an efficient implementation of a `DataFrame`. `DataFrame`s are essentially multidimensional arrays with attached row and column labels, and often with heterogeneous types and/or missing data. As well as offering a convenient storage interface for labeled data, Pandas implements a number of powerful data operations familiar to users of both database frameworks and spreadsheet programs.

As we saw, NumPy's `ndarray` data structure provides essential features for the type of clean, well-organized data typically seen in numerical computing tasks. While it serves this purpose very well, its limitations become clear when we need more flexibility (e.g. attaching labels to data, working with missing data, etc.), each of which is an important piece of analyzing the less structured data available in many forms in the world around us. Pandas, and in particular its `Series` and `DataFrame` objects, builds on the NumPy array structure and provides efficient access to these sorts of "data munging" tasks that occupy much of a data scientist's time.

In this notebook, we will focus on the mechanics of using `Series`, `DataFrame`, and related structures effectively. We will use examples drawn from real datasets where appropriate, but these examples are not necessarily the focus.

Once Pandas is installed, you can import it along with NumPy and check the version. Just as we generally import NumPy under the alias `np`, we will import Pandas under the alias `pd`. This import convention will be used throughout the remainder of this course.

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

```
Out[1]: '1.0.5'
```

The Pandas Series Object

A Pandas `Series` is a one-dimensional array of indexed data. It can be created from a list or array as follows:

```
In [ ]: data = pd.Series([0.25, 0.5, 0.75, 1.0],
data)
```

As we see in the output, the `Series` wraps both a sequence of values and a sequence of indices, which we can access with the `values` and `index` attributes. The values are simply a familiar NumPy array. The `index` is an array-like object of type `pd.Index`, which we'll discuss in more detail momentarily. Like with a NumPy array, data can be accessed by the associated index via the familiar Python square-bracket notation.

```
In [ ]: print(data.values)
print(data.index)
print(data['b'])
```

Series as generalized NumPy array

From what we've seen so far, it may look like the `Series` object is basically interchangeable with a one-dimensional NumPy array. The essential difference is the presence of the index: while the NumPy Array has an *implicitly defined* integer index used to access the values, the `Pandas Series` has an *explicitly defined* index associated with the values.

This explicit index definition gives the `Series` object additional capabilities. For example, the index need not be an integer, but can consist of values of any desired type. For example, if we wish, we can use strings as an index:

```
In [ ]: data = pd.Series([0.25, 0.5, 0.75, 1.0],
index=['a', 'b', 'c', 'd'],
data)
```

```
print(data)
print(data['b'])
```

Series as specialized dictionary

In this way, you can think of a `Pandas Series` a bit like a specialization of a Python dictionary. A dictionary is a structure that maps arbitrary keys to a set of arbitrary values, and a `Series` is a structure which maps typed keys to a set of typed values. This typing is important: just as the type-specific compiled code behind a NumPy array makes it more efficient than a Python list for certain operations, the type information of a `Pandas Series` makes it much more efficient than Python dictionaries for certain operations.

The `Series`-as-dictionary analogy can be made even more clear by constructing a `Series` object directly from a Python dictionary:

```
In [ ]: population_dict = {'California': 38332521,
'Texas': 26448193,
'New York': 19651127,
'Florida': 19552860,
'Illinois': 12882135}
population = pd.Series(population_dict)
```

By default, a `Series` will be created where the index is drawn from the sorted keys. From here, typical dictionary-style item access can be performed. Unlike a dictionary, though, the `Series` also supports array-style operations such as slicing:

```
In [ ]: population['California':'New York']
```

The Pandas DataFrame Object

The next fundamental structure in Pandas is the `DataFrame`. Like the `Series` object discussed in the previous section, the `DataFrame` can be thought of either as a generalization of a NumPy array, or as a specialization of a Python dictionary. We'll now take a look at each of these perspectives.

DataFrame as a generalized NumPy array

If a `Series` is an analog of a one-dimensional array with flexible indices, a `DataFrame` is an analog of a two-dimensional array with both flexible row indices and flexible column names. Just as you might think of a two-dimensional array as an ordered sequence of aligned one-dimensional columns, you can think of a `DataFrame` as a sequence of aligned `Series` objects. Here, by "aligned" we mean that they share the same index.

To demonstrate this, let's first construct a new `Series` listing the area of each of the five states discussed in the previous section:

```
In [ ]: area_dict = {'California': 423967, 'Texas': 695662, 'New York': 141297,
'Florida': 170312, 'Illinois': 149995}
area = pd.Series(area_dict)
area
```

Now that we have this along with the `population` `Series` from before, we can use a dictionary to construct a single two-dimensional object containing this information:

```
In [ ]: states = pd.DataFrame({'population': population,
'area': area})
states
```

Like the `Series` object, the `DataFrame` has an `index` attribute that gives access to the index labels. Additionally, the `DataFrame` has a `columns` attribute, which is an `Index` object holding the column labels. Thus the `DataFrame` can be thought of as a generalization of a two-dimensional NumPy array, where both the rows and columns have a generalized index for accessing the data.

```
In [ ]: print(states.index)
print(states.columns)
```

The Pandas Index Object

We have seen here that both the `Series` and `DataFrame` objects contain an explicit *index* that lets you reference and modify data. This `Index` object is an interesting structure in itself, and it can be thought of either as an *immutable array* or as an *ordered set* (technically a multi-set, as `Index` objects may contain repeated values). Those views have some interesting consequences in the operations available on `Index` objects. As a simple example, let's construct an `Index` from a list of integers:

```
In [ ]: ind = pd.Index([2, 3, 5, 7, 11])
ind
```

```
Out[ ]: 2
3
5
7
11
```

Data Indexing and Selection

In the previous notebook, we looked in detail at methods and tools to access, set, and modify values in NumPy arrays. Here we'll look at similar means of accessing and modifying values in Pandas `Series` and `DataFrame` objects. If you have used the NumPy patterns, the corresponding patterns in Pandas will feel very familiar, though there are a few quirks to be aware of.

We'll start with the simple case of the one-dimensional `Series` object, and then move on to the more complicated two-dimensional `DataFrame` object.

Data Selection in Series

As we saw in the previous section, a `Series` object acts in many ways like a one-dimensional NumPy array, and in many ways like a standard Python dictionary. If we keep these two overlapping analogies in mind, it will help us to understand the patterns of data indexing and selection in these arrays.

Series as dictionary

Like a dictionary, the `Series` object provides a mapping from a collection of keys to a collection of values:

```
In [ ]: import pandas as pd
data = pd.Series([0.25, 0.5, 0.75, 1.0],
index=['a', 'b', 'c', 'd'])
data
```

```
In [ ]: # We can also use dictionary-like Python expressions and methods to examine the keys/indices and values
print(data.keys())
print(list(data.items()))
```

```
In [ ]: data['e'] = 1.25
data
```

Indexers: loc and iloc

The slicing and indexing conventions can be a source of confusion. For example, if your `Series` has an explicit integer index, an indexing operation such as `data[1]` will use the explicit indices, while a slicing operation like `data[1:3]` will use the implicit Python-style index.

```
In [ ]: data = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])
data
```

```
In [ ]: # explicit index when indexing
data[1]
```

```
In [ ]: # implicit index when slicing
data[1:3]
```

Because of this potential confusion in the case of integer indexes, Pandas provides some special *indexer* attributes that explicitly expose certain indexing schemes. These are not functional methods, but attributes that expose a particular slicing interface to the data in the `Series`.

First, the `loc` attribute allows indexing and slicing that always references the explicit index:

```
In [ ]: data.loc[1]
In [ ]: data.loc[1:3]
```

The `iloc` attribute allows indexing and slicing that always references the implicit Python-style index:

```
In [ ]: data.iloc[1]
In [ ]: data.iloc[1:3]
```

One guiding principle of Python code is that "explicit is better than implicit." The explicit nature of `loc` and `iloc` make them very useful in maintaining clean and readable code; especially in the case of integer indexes, I recommend using these both to make code easier to read and understand, and to prevent subtle bugs due to the mixed indexing/slicing convention.

Data Selection in DataFrame

Recall that a `DataFrame` acts in many ways like a two-dimensional or structured array, and in other ways like a dictionary of `Series` structures sharing the same index. These analogies can be helpful to keep in mind as we explore data selection within this structure.

DataFrame as a dictionary

The first analogy we will consider is the `DataFrame` as a dictionary of related `Series` objects. Let's return to our example of areas and populations of states:

```
In [ ]: area = pd.Series({'California': 423967, 'Texas': 695662,
'New York': 141297, 'Florida': 170312,
'Illinois': 149995})
pop = pd.Series({'California': 38332521, 'Texas': 26448193,
'New York': 19651127, 'Florida': 19552860,
'Illinois': 12882135})
data = pd.DataFrame({'area':area, 'pop':pop})
data
```

The individual `Series` that make up the columns of the `DataFrame` can be accessed via dictionary-style indexing of the column name:

```
In [ ]: data['area']
```

Equivalently, we can use attribute-style access with column names that are strings:

```
In [ ]: data.area
```

DataFrame as two-dimensional array

As mentioned previously, we can also view the `DataFrame` as an enhanced two-dimensional array. We can examine the raw underlying data array using the `values` attribute:

```
In [ ]: data.values
```

For array-style indexing, we can use the `loc`, `iloc`, and `ix` indexers mentioned earlier. By using the `iloc` indexer, we can index the underlying array as if it is a simple NumPy array (using the implicit Python-style index), but the `DataFrame` index and column labels are maintained in the result:

```
In [ ]: data.iloc[1:, 1:]
```

Similarly, using the `loc` indexer we can index the underlying data in an array-like style but using the explicit index and column names:

```
In [ ]: data.loc['Illinois', : 'pop']
```

Any of the familiar NumPy-style data access patterns can be used within these indexers. For example, in the `loc` indexer we can combine masking and fancy indexing as in the following:

```
In [ ]: data.loc[data.density > 100, ['pop', 'density']]
```

Any of these indexing conventions may also be used to set or modify values; this is done in the standard way that you might be accustomed to from working with NumPy:

```
In [ ]: data.iloc[0, 2] = 90
data
```

To build up your fluency in Pandas data manipulation, I suggest spending some time with a simple `DataFrame` and exploring the types of indexing, slicing, masking, and fancy indexing that are allowed by these various indexing approaches.

```
Out[ ]:
```

Operating on Pieces in Pandas

One of the essential pieces of NumPy is the ability to perform quick element-wise operations, both with basic arithmetic (addition, subtraction, multiplication, etc.) and with more sophisticated operations (trigonometric functions, exponential and logarithmic functions, etc.). Pandas inherits much of this functionality from NumPy, and the `ufuncs` that we introduced in NumPy notebook are key to this.

Pandas includes a couple useful twists, however: for unary operations like negation function, these `ufuncs` will *preserve index and column labels* in the output, and for binary operations such as addition and multiplication, Pandas will automatically *align indices* when passing the objects to the `ufunc`. This means that keeping the context of data and combining data from different sources—both potentially error-prone tasks with raw NumPy arrays—become essentially foolproof ones with Pandas. We will additionally see that there are well-defined operations between one-dimensional `Series` structures and two-dimensional `DataFrame` structures.

Ufuncs: Index Preservation

Because Pandas is designed to work with NumPy, any NumPy `ufunc` will work on Pandas `Series` and `DataFrame` objects. Let's demonstrate this in a `Series`:

```
In [ ]: rng = np.random.RandomState(42) # to reproduce same output
ser = pd.Series(rng.randint(0, 10, 4))
ser
```

If we apply a NumPy `ufunc` on either of these objects, the result will be another Pandas object with the *indices preserved*:

```
In [ ]: np.exp(ser) # exponential function, calculates e^x
```

UFuncs: Index Alignment

For binary operations on two `Series` or `DataFrame` objects, Pandas will align indices in the process of performing the operation. This is a very convenient when working with incomplete data, as we'll see in some of the examples that follow.

Index alignment in Series

As an example, suppose we are combining two different data sources, and find only the top three US states by area and the top three US states by population:

```
In [ ]: area = pd.Series({'Alaska': 1723337, 'Texas': 695662,
'California': 423967}, name='area')
population = pd.Series({'California': 38332521, 'Texas': 26448193,
'New York': 19651127}, name='population')
```

When we divide these to compute the population density, the resulting array contains the union of indices of the two input arrays:

```
In [ ]: population / area
```

Index alignment in DataFrame

A similar type of alignment takes place for both columns and indices when performing operations on `DataFrame`s:

```
In [ ]: A = pd.DataFrame(rng.randint(0, 20, (2, 2)),
columns=list('AB'))
A
```

```
In [ ]: B = pd.DataFrame(rng.randint(0, 10, (3, 3)),
columns=list('BAC'))
B
```

```
In [ ]: A + B
```

Notice that indices are aligned correctly irrespective of their order in the two objects, and indices in the result are sorted. We can use the associated object's `arithmetic` method and pass any desired `fill_value` to be used in place of missing entries. Here we'll fill with the mean of all values in `A` (computed by first stacking the rows of `A`):

```
In [ ]: fill = A.stack().mean()
A.add(B, fill_value=fill)
```

Ufuncs: Operations Between DataFrame and Series

When performing operations between a `DataFrame` and a `Series`, the index and column alignment is similarly maintained. Operations between a `DataFrame` and a `Series` are similar to operations between a two-dimensional and one-dimensional NumPy array. Consider one common operation, where we find the difference of a two-dimensional array and one of its rows:

```
In [ ]: A = rng.randint(10, size=(3, 4))
A
```

```
In [ ]: A - A[0]
```

In Pandas, the convention similarly operates row-wise by default:

```
In [ ]: df = pd.DataFrame(A, columns=list('QRST'))
df - df.iloc[0]
```

If you would instead like to operate column-wise, you can use the object methods mentioned earlier, while specifying the `axis` keyword:

```
In [ ]: df.subtract(df['R'], axis=0)
```

This presentation and alignment of indices and columns means that operations on data in Pandas will always maintain the data context, which prevents the types of slip errors that might come up when working with heterogeneous and/or misaligned data in raw NumPy arrays.

Handling Missing Data

The difference between data found in many tutorials and data in the real world is that real-world data is rarely clean and homogeneous. In particular, many interesting datasets will have some amount of data missing. To make matters even more complicated, different data sources may indicate missing data in different ways.

In this section, we will discuss some general considerations for missing data, discuss how Pandas chooses to represent it, and demonstrate some built-in Pandas tools for handling missing data in Python. Here and throughout the course, we'll refer to missing data in general as *null*, *NaN*, or *NA* values.

Missing Data in Pandas

The way in which Pandas handles missing values is constrained by its reliance on the NumPy package, which does not have a built-in notion of NA values for non-floating-point data types.

Pandas could have followed R's lead in specifying bit patterns for each individual data type to indicate nullness, but this approach turns out to be rather unwieldy. While R contains four basic data types, NumPy supports far more than this: for example, while R has a single integer type, NumPy supports fourteen basic integer types once you account for available precisions, signedness, and endianness of the encoding. Reserving a specific bit pattern in all available NumPy types would lead to an unwieldy amount of overhead in special-casing various operations for various types, likely even requiring a new fork of the NumPy package. Further, for the smaller data types (such as 8-bit integers), sacrificing a bit to use as a mask would significantly reduce the range of values it can represent.

NumPy does have support for masked arrays—that is, arrays that have a separate Boolean mask array attached for marking data as "good" or "bad." Pandas could have derived from this, but the overhead in both storage, computation, and code maintenance makes that an unattractive choice.

With these constraints in mind, Pandas chose to use sentinels for missing data, and this choice has some side effects, as we will see, but in practice ends up being a good compromise in most cases of interest.

Operating on Pandas Null Values

As we have seen, Pandas treats `None` and `NaN` as essentially interchangeable for indicating missing or null values. To facilitate this convention, there are several useful methods for detecting, removing, and replacing null values in Pandas data structures. They are:

- `isnull()`: Generate a boolean mask indicating missing values
- `notnull()`: Opposite of `isnull()`
- `dropna()`: Return a filtered version of the data
- `fillna()`: Return a copy of the data with missing values filled or imputed

We will conclude this section with a brief exploration and demonstration of these routines.

Detecting null values

Pandas data structures have two useful methods for detecting null data: `isnull()` and `notnull()`. Either one will return a Boolean mask over the data. For example:

```
In [ ]: data = pd.Series([1, np.nan, 'hello', None])
In [ ]: data.isnull()
In [ ]: data.notnull()
```

The `isnull()` and `notnull()` methods produce similar Boolean results for `DataFrame`s.

Dropping null values

In addition to the masking used before, there are the convenience methods, `dropna()` (which removes NA values) and `fillna()` (which fills in NA values). For a `Series`, the result is straightforward:

```
In [ ]: data.dropna()
```

For a `DataFrame`, there are more options. Consider the following `DataFrame`:

```
In [ ]: df = pd.DataFrame([[1, np.nan, 2],
[2, 3, 5],
[3, np.nan, 4, 6]])
df
```

We cannot drop single values from a `DataFrame`; we can only drop full rows or full columns. Depending on the application, you might want one or the other, so `dropna()` gives a number of options for a `DataFrame`:

```
In [ ]: df.dropna() # By default, dropna() will drop all rows in which any null value is present
In [ ]: df.dropna(axis='columns') # Alternatively, you can drop NA values along all columns containing a null value
```

But this drops some good data as well; you might rather be interested in dropping rows or columns with *all* NA values, or a majority of NA values. This can be specified through the `how` or `thresh` parameters, which allow fine control of the number of nulls to allow through. Try them out on your own.

Filling null values

Sometimes rather than dropping NA values, you'd rather replace them with a valid value. This value might be a single number like zero, or it might be some sort of interpolation or extrapolation from the good values. You could do this in-place using the `isnull()` method as a mask, but because it is such a common operation Pandas provides the `fillna()` method, which returns a copy of the array with the null values replaced.

Consider the following `Series`:

```
In [ ]: data = pd.Series([1, np.nan, 2, None, 3], index=list('abcde'))
data
```

We can fill NA entries with a single value, such as zero:

```
In [ ]: data.fillna(0)
```

For `DataFrame`s, the options are similar, but we can also specify an `axis` along which the fills take place:

```
In [ ]: df
In [ ]: df.fillna(method='ffill', axis=1)
```

Notice that if a previous value is not available during a forward fill, the NA value remains.

Combining Datasets: Concat and Append

Some of the most interesting studies of data come from combining different data sources. These operations can involve anything from very straightforward concatenation of two different datasets, to more complicated database-style joins and merges that handle any overlaps between the datasets. `Series` and `DataFrame`s are built with this type of operation in mind, and Pandas includes functions and methods that make this sort of data wrangling fast and straightforward.

Here we'll take a look at simple concatenation of `Series` and `DataFrame`s with the `pd.concat` function; later we'll dive into more sophisticated in-memory merges and joins implemented in Pandas.

For convenience, we'll define this function which creates a `DataFrame` of a particular form that will be useful below:

```
In [2]: def make_df(cols, ind):
"""Quickly make a DataFrame"""
data = {}
for c in cols:
data[c] = str(c) + str(ind)
return pd.DataFrame(data, ind)
# example DataFrame, range(3))
```

```
Out[2]:
```

	A	B	C
0	A0	B0	C0
1	A1	B1	C1
2	A2	B2	C2

In addition, we'll create a quick class that allows us to display multiple `DataFrame`s side by side. The code makes use of the special `repr_html` method, which IPython uses to implement its rich object display. The use of this will become clearer as we continue our discussion in the following section.

```
In [3]: class display(object):
"""Display HTML representation of multiple objects"""
template = """<div>{font-family:'Courier New', Courier, monospace}>{0}</div>{1}</div>"""
def __init__(self, args):
self.args = args
def _repr_html_(self):
return "\n".join(self.template.format(a, eval(a)._repr_html_())
for a in self.args)
def _repr_text_(self):
return "\n\n".join(a + '\n' + repr(eval(a))
for a in self.args)
```

Simple Concatenation with pd.concat

Pandas has data structures, `pd.concat()`, which has a similar syntax to `np.concatenate`, but contains a number of options. `pd.concat()` can be used for a simple concatenation of `Series` or `DataFrame` objects, just as `np.concatenate()` can be used for simple concatenations of arrays:

```
In [9]: df1 = make_df('AB', [1, 2])
df2 = make_df('BCD', [3, 4])
display(df1, df2, "pd.concat(df1, df2)")
```

```
Out[9]:
```

	A	B		C	D
df1			pd.concat(df1, df2)		
	A	B		A	B
1	A1	B1	3	A3	B3
2	A2	B2	4	A4	B4
			3	A3	B3
			4	A4	B4

By default, the concatenation takes place row-wise within the `DataFrame` (i.e. `axis=0`). Like `np.concatenate`, `pd.concat` allows specification of an axis along which concatenation will take place. Consider the following example:

```
In [7]: df3 = make_df('AB', [0, 1])
df4 = make_df('CD', [0, 1])
display(df3, df4, "pd.concat(df3, df4, axis=1)")
```

```
Out[7]:
```

	A	B		C	D
df3			pd.concat(df3, df4, axis=1)		
	A	B		A	B
0	A0	B0	0	C0	D0
1	A1	B1	1	C1	D1

Ignoring the index

Sometimes the index itself does not matter, and you would prefer it to simply be ignored. This option can be specified using the `ignore_index` flag. With this set to true, the concatenation will create a new integer index for the resulting `Series`.

```
In [11]: x = make_df('AB', [0, 1])
y = make_df('AB', [2, 3])
display(x, y, "pd.concat(x, y, ignore_index=True)")
```



```
In [18]: df4 = pd.DataFrame({'group': ['Accounting', 'Engineering', 'HR'],
display('df3', 'df4', 'pd.merge(df3, df4)')
```

Out [18]:

df3				df4							
employee	group	hire_date		group	supervisor	employee	group	hire_date	supervisor		
0	Bob	Accounting	2008	0	Accounting	Carly	0	Bob	Accounting	2008	Carly
1	Jake	Engineering	2012	1	Engineering	Guido	1	Jake	Engineering	2012	Guido
2	Lisa	Engineering	2004	2	HR	Steve	2	Lisa	Engineering	2004	Guido
3	Sue	HR	2014				3	Sue	HR	2014	Steve

The resulting `DataFrame` has an additional column with the "supervisor" information, where the information is repeated in one or more locations as required by the inputs.

Many-to-many joins

The resulting `DataFrame` has an additional column with the "supervisor" information, where the information is repeated in one or more locations as required by the inputs.

Many-to-many joins

Many-to-many joins are a bit confusing conceptually, but are nevertheless well defined. If the key column in both the left and right array contains duplicates, then the result is a many-to-many merge. This will be perhaps most clear with a concrete example. Consider the following, where we have a `DataFrame` showing one or more skills associated with a particular group. By performing a many-to-many join, we can recover the skills associated with any individual person:

```
In [19]: df5 = pd.DataFrame({'group': ['Accounting', 'Accounting',
display('df1', 'df5', 'pd.merge(df1, df5)')
```

Most simply, you can explicitly specify the name of the key column using the `on` keyword, which takes a column name or a list of column names. This option works only for both the left and right `DataFrame`s that have the specified column name.

```
In [20]: display('df1', 'df2', 'pd.merge(df1, df2, on='employee')')
```

```
Out[20]:
```

df1			df2			pd.merge(df1, df2, on='employee')		
employee	group	hire_date	employee	group	hire_date	employee	group	hire_date
0	Bob	Accounting	0	Lisa	2004	0	Bob	Accounting
1	Jake	Engineering	1	Bob	2008	1	Jake	Engineering
2	Lisa	Engineering	2	Jake	2012	2	Lisa	Engineering
3	Sue	HR	3	Sue	2014	3	Sue	HR

Specifying Set Arithmetic for Joins

In all the preceding examples we have glossed over one important consideration in performing a join: the type of set arithmetic used in the join. This comes up when a value appears in one key column but not the other.

Consider the below example, where we merge two datasets that have only a single "name" entry in common. Mary. By default, the result contains the *intersection* of the two sets of inputs; this is what is known as an *inner* join. We can specify this explicitly using the `how` keyword, which defaults to `"inner"`:

```
In [22]: df6 = pd.DataFrame({'name': ['Peter', 'Paul', 'Mary'],
                             'food': ['fish', 'beans', 'bread']},
                             columns=['name', 'food'])
df7 = pd.DataFrame({'name': ['Mary', 'Joseph'],
                    'drink': ['wine', 'beer']},
                    columns=['name', 'drink'])
```

These three types of joins can be used with other Pandas tools to implement a wide array of functionality. But in practice, datasets are rarely as clean as the one we're working with here.

The `on` keyword

Most simply, you can explicitly specify the name of the key column using the `on` keyword, which takes a column name or a list of column names. This option works only if both the left and right `DataFrame`s have the specified column name.

```
In [20]: display('df1', 'df2', "pd.merge(df1, df2, on='employee')")
```

2 Mary bread 3 Mary wine

3 Joseph NaN beer

The `left join` and `right join` return joins over the left entries and right entries, respectively. In the below example, the output rows now correspond to the entries in the left input. Using `how='right'` works in a similar manner.

```
In [16]: display(df6, 'df7', 'pd.merge(df6, df7, how='left')')
```

Out[16]:

df6			df7			pd.merge(df6, df7, how='left')			
name	food		name	drink		name	food	drink	
0	Peter	fish	0	Mary	wine	0	Peter	fish	NaN
1	Paul	beans	1	Joseph	beer	1	Paul	beans	NaN
2	Mary	bread				2	Mary	bread	wine

Specifying Set Arithmetic for Joins

In all the preceding examples we have glossed over one important consideration in performing a join: the type of set arithmetic used in the join. This comes up when a value appears in one key column but not the other.

Consider the below example, where we merge two datasets that have only a single "name" entry in common: Mary. By default, the result contains the *intersection* of the two sets of inputs; this is what is known as an *inner join*. We can specify this explicitly by the `how` keyword, which defaults to "inner":

```
In [22]: df6 = pd.DataFrame({'name': ['Peter', 'Paul', 'Mary'],
df7 = pd.DataFrame({'name': ['Mary', 'Joseph'],
display('df6', 'df7', 'pd.merge(df6, df7, how='inner')')
```

Out [22]:

df6			df7			pd.merge(df6, df7, how='inner')			
name	food		name	drink		name	food	drink	
0	Peter	fish	0	Mary	wine	0	Mary	bread	wine
1	Paul	beans	1	Joseph	beer				
2	Mary	bread							

Other options for the `how` keyword are "outer", "left", and "right". An *outer join* returns a join over the union of the input columns, and fills in all missing values with NAs:

```
In [15]: display('df6', 'df7', "pd.merge(df6, df7, how='outer')")
```

2	AL	total	2012	4817528.0	Alabama
3	AL	under18	2010	1130986.0	Alabama
4	AL	total	2010	4785570.0	Alabama
5	AL	under18	2011	1125763.0	Alabama

Let's double-check whether there were any mismatches here, which we can do by looking for rows with nulls:

```
In [25]: merged.isnull().any()
```

```
Out[25]: state/region    False
         ages           False
         year            False
         population      True
         state           True
         dtype: bool
```

The *left join* and *right join* return joins over the left entries and right entries, respectively. In the below example, the output rows now correspond to the entries in the left input. Using `how='right'` works in a similar manner.

```
In [16]: display('df6', 'df7', "pd.merge(df6, df7, how='left')")
```

Out [16]:

df6			df7			pd.merge(df6, df7, how='left')			
name	food		name	drink		name	food	drink	
0	Peter	fish	0	Mary	wine	0	Peter	fish	NaN
1	Paul	beans	1	Joseph	beer	1	Paul	beans	NaN
2	Mary	bread	2	Mary	bread	wine			
			3	Joseph	NaN	beer			

All of these options can be applied straightforwardly to any of the preceding join types.

Example: US States Data

Merge and join operations come up most often when combining data from different sources. Here we will consider an example of some data about US states and their populations. The CSV files can be found under the data folder you have downloaded along with the notebooks.

Let's take a look at the three datasets, using the Pandas `read_csv()` function:

```
In [23]: pop = pd.read_csv('data/state-population.csv')
areas = pd.read_csv('data/state-areas.csv')
abbrevs = pd.read_csv('data/state-abbrevs.csv')
display('pop.head()', 'areas.head()', 'abbrevs.head()') # by default, head method displays the first five rows
```

0	AL	under18	2012	1117489.0	Alabama	52423.0
1	AL	total	2012	4817528.0	Alabama	52423.0
2	AL	under18	2010	1130966.0	Alabama	52423.0
3	AL	total	2010	4785570.0	Alabama	52423.0
4	AL	under18	2011	1125763.0	Alabama	52423.0

Again, let's check for nulls to see if there were any mismatches:

```
In [30]: final.isnull().any()
```

```
Out[30]: state/region    False
ages                  False
year                  False
population             True
state                  False
area (sq. mi)         True
dtype: bool
```

There are nulls in the `area` column, we can take a look to see which regions were ignored here:

Given this information, say we want to compute a relatively straightforward result: rank US states and territories by their 2010 population density. We clearly have the data here to find this result, but we'll have to combine the datasets to find the result.

We'll start with a many-to-one merge that will give us the full state name within the population `DataFrame`. We want to merge based on the `state/region` column of `pop`, and the `abbreviation` column of `abbrevs`. We'll use `how='outer'` to make sure no data is thrown away due to mismatched labels.

```
In [24]: merged = pd.merge(pop, abbrevs, how='outer',
merged = merged.drop('abbreviation', 1) # drop duplicate info
merged.head()
```

Now we have all the data we need. To answer the question of interest, let's first select the portion of the data corresponding with the year 2000, and the total population. We'll use the `query()` function to do this quickly.

```
In [37]: data2010 = final.query("year == 2010 & ages == 'total'")
data2010.head()
```

Out[37]:

	state/region	ages	year	population	state	area (sq. mi)
3	AL	total	2010	4785570.0	Alabama	52423.0
91	AK	total	2010	713868.0	Alaska	666425.0
101	AZ	total	2010	6408790.0	Arizona	114006.0
189	AR	total	2010	2922280.0	Arkansas	53182.0
197	CA	total	2010	37338601.0	California	163707.0

Let's double-check whether there were any mismatches here, which we can do by looking for rows with nulls:

```
In [25]: merged.isnull().any()
```

Out [25]:

state/region	False
ages	False
year	False
population	True
state	True
dtype:	bool

Some of the `population` info is null; let's figure out which these are!

```
In [26]: merged[merged['population'].isnull()].head()
```

We can also check the end of the list:

```

In [36]: density.tail()

Out[36]:
state
South Dakota    10.583512
North Dakota    9.537565
Montana         6.736171
Wyoming         5.768079
Alaska          1.087509
dtype: float64

```

We see that the least dense state, by far, is Alaska, averaging slightly over one resident per square mile.

This type of messy data merging is a common task when trying to answer questions using real-world data sources. I hope that this example has given you an idea of the ways you can combine tools we've covered in order to gain insight from your data!

It appears that all the null population values are from Puerto Rico prior to the year 2000; this is likely due to this data not being available from the original source.

More importantly, we see also that some of the new `state` entries are also null, which means that there was no corresponding entry in the `abbrevs` key! Let's figure out which regions lack this match:

```
In [27]: merged.loc[merged['state'].isnull()]['state'].unique()
```

Out [27]:

array(['PR', 'USA'], dtype=object)

We can quickly infer the issue: our population data includes entries for Puerto Rico (PR) and the United States as a whole (USA), while these entries do not appear in the state abbreviation key. We can fix these quickly by filling in appropriate entries:

```
In [28]: merged.loc[merged['state/region'] == 'PR', 'state'] = 'Puerto Rico'
merged.loc[merged['state/region'] == 'USA', 'state'] = 'United States'
merged.isnull().any()
```

Out [28]:

state/region	False
ages	False
year	False
population	True
state	False
dtype:	bool

No more nulls in the `state` column: we're all set!

Now we can merge the result with the area data using a similar procedure. Examining our results, we will want to join on the `state` column in both:

```
In [29]: final = pd.merge(merged, areas, on='state', how='left')
final.head()
```

```
Out[5]: 2.8119254917081569
```

```
In [6]: ser.mean()
```

```
Out[6]: 0.56238509834163142
```

For a `DataFrame`, by default the aggregates return results within each column:

```
In [7]: df = pd.DataFrame({'A': rng.rand(5),
                           'B': rng.rand(5)})
df
```

```
Out[7]:
```

	A	B
0	0.155995	0.020584
1	0.058084	0.969910

Again, let's check for nulls to see if there were any mismatches:

```
In [30]: final.isnull().any()
```

Out [30]:

state/region	False
ages	False
year	False
population	True
state	False
area (sq. mi)	True
dtype:	bool

There are nulls in the `area` column; we can take a look to see which regions were ignored here:

```
In [31]: final['state'][final['area (sq. mi)'].isnull()].unique()
```

Out [31]:

array(['United States'], dtype=object)
--

We see that our `areas` `DataFrame` does not contain the area of the United States as a whole. We could insert the appropriate value (using the sum of all state areas, for instance), but in this case we'll just drop the null values because the population density of the entire United States is not relevant to our current discussion:

```
In [32]: final.dropna(inplace=True)
final.head()
```

Out [32]:

state/region	ages	year	population	state	area (sq. mi)	
0	AL	under18	2012	1117489.0	Alabama	52423.0
1	AL	total	2012	4817528.0	Alabama	52423.0
2	AL	under18	2010	1130966.0	Alabama	52423.0
3	AL	total	2010	4785570.0	Alabama	52423.0
4	AL	under18	2011	1125763.0	Alabama	52423.0

Now we have all the data we need. To answer the question of interest, let's first select the portion of the data corresponding with the year 2000, and the total population. We'll use the `query()` function to do this quickly.

```
In [37]: data2010 = final.query("year == 2010 & ages == 'total'")
data2010.head()
```

This is a good time to stop and dig into the above discussed Pandas concepts, try out the variations in each of these code snippets and evaluate the individual steps to make sure you understand exactly what they are doing to the result. There are several other advanced topics in Pandas, which you can slowly explore once you master the basics.

This notebook contains an excerpt from the [Python Data Science Handbook](#) by Jake VanderPlas. The text is released under the [CC-BY-NC-ND license](#), and code is released under the [MIT license](#).

Now let's compute the population density and display it in order. We'll start by re-indexing our data on the state, and then compute the result:

```
In [34]: data2010.set_index('state', inplace=True)
density = data2010['population'] / data2010['area (sq. mi)']
```

```
In [35]: density.sort_values(ascending=False, inplace=True)
density.head()
```

Out [35]:

state	
District of Columbia	8898.897059
Puerto Rico	1058.665149
New Jersey	1009.233268
Rhode Island	651.339159
Connecticut	645.600649
dtype:	float64

The result is a ranking of US states plus Washington, DC, and Puerto Rico in order of their 2010 population density, in residents per square mile. We can see that by far the densest region in this dataset is Washington, DC (i.e., the District of Columbia); among states, the densest is New Jersey.

We can also check the end of the list:

```
In [36]: density.tail()
```

Out [36]:

state	
South Dakota	10.583512
North Dakota	9.537565
Montana	6.736171
Wyoming	5.746079
Alaska	1.087509
dtype:	float64

We see that the least dense state, by far, is Alaska, averaging slightly over one resident per square mile.

This type of messy data merging is a common task when trying to answer questions using real-world data sources. I hope that this example has given you an idea of the ways you can combine tools we've covered in order to gain insight from your data!

Aggregation and Grouping

An essential piece of analysis of large data is efficient summarization: computing aggregations like `sum()`, `mean()`, `median()`, `min()`, and `max()`, in which a single number gives insight into the nature of a potentially large dataset. In this section, we'll explore aggregations in Pandas.

Planets Data

Here we will use the Planets dataset, available via the [Seaborn](#) package. It gives information on planets that astronomers have discovered around other stars (known as *extrasolar planets* or *exoplanets* for short). It can be downloaded with a simple Seaborn command:

```
In [38]: import seaborn as sns
planets = sns.load_dataset('planets')
planets.shape
```

Out [38]:

(1035, 6)

```
Out [39]: planets.head() # This has some details on the 1,000+ extrasolar planets discovered up to 2014.
```

Out [39]:

method	number	orbital_period	mass	distance	year	
0	Radial Velocity	1	269.300	7.10	77.40	2006
1	Radial Velocity	1	874.774	2.21	58.95	2008
2	Radial Velocity	1	763.000	2.60	19.84	2011
3	Radial Velocity	1	326.030	19.40	110.62	2007
4	Radial Velocity	1	516.220	10.50	119.47	2009

Simple Aggregation in Pandas

Earlier, we explored some of the data aggregations available for NumPy arrays. As with a one-dimensional NumPy array, for a Pandas `Series`, the aggregates return a single value:

```
In [4]: rng = np.random.RandomState(42)
ser = pd.Series(rng.rand(5))
ser
```

Out [4]:

0	0.374540
1	0.959714
2	0.731994
3	0.598658
4	0.156019
dtype:	float64

Out [5]:

ser.sum()

Out [6]:

ser.mean()

Out [6]:

0.56238509834163142

For a `DataFrame`, by default the aggregates return results within each column:

```
In [7]: df = pd.DataFrame({'A': rng.rand(5),
df
```

Out [7]:

	A	B
0	0.155995	0.020584
1	0.058084	0.969910
2	0.866176	0.832443
3	0.601115	0.212339
4	0.708073	0.181825

Out [8]:

df.mean()

Out [8]:

A	0.477888
B	0.443420
dtype:	float64

By specifying the `axis` argument, you can instead aggregate within each row:

```
In [9]: df.mean(axis='columns')
```

Out [9]:

0	0.088290
1	0.513987
2	0.849309
3	0.406727
4	0.444949
dtype:	float64

Pandas `Series` and `DataFrame`s include all of the common aggregates mentioned in NumPy notebook; in addition, there is a convenience method `describe()` that computes several common aggregates for each column and returns the result. Let's use this on the Planets data, for now dropping rows with missing values: