

Data Indexing and Selection

In this notebook we will look at means of accessing and modifying values in Pandas `Series` and `DataFrame` objects.

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

```
Out[1]: a    0.25
        b    0.50
        c    0.75
        d    1.00
        dtype: float64
```

```
In [2]: data['b']
```

```
Out[2]: 0.5
```

We can also use dictionary-like Python expressions and methods to examine the keys/indices and values:

```
In [3]: 'a' in data
```

```
Out[3]: True
```

```
In [4]: data.keys()
```

```
Out[4]: Index(['a', 'b', 'c', 'd'], dtype='object')
```

```
In [5]: list(data.items())
```

```
Out[5]: [('a', 0.25), ('b', 0.5), ('c', 0.75), ('d', 1.0)]
```

`Series` objects can even be modified with a dictionary-like syntax. Just as you can extend a dictionary by assigning to a new key, you can extend a `Series` by assigning to a new index value:

```
In [6]: data['e'] = 1.25
data
```

```
Out[6]: a    0.25
        b    0.50
        c    0.75
        d    1.00
        e    1.25
        dtype: float64
```

This easy mutability of the objects is a convenient feature: under the hood, Pandas is making decisions about memory layout and data copying that might need to take place; the user generally does not need to worry about these issues.

Series as one-dimensional array

A `Series` builds on this dictionary-like interface and provides array-style item selection via the same basic mechanisms as NumPy arrays – that is, *slices*, *masking*, and *fancy indexing*. Examples of these are as follows:

```
In [7]: # slicing by explicit index
data['a':'c']
```

```
Out[7]: a    0.25
        b    0.50
        c    0.75
        dtype: float64
```

```
In [8]: # slicing by implicit integer index
data[0:2]
```

```
Out[8]: a    0.25
        b    0.50
        dtype: float64
```

```
In [9]: # masking
data[(data > 0.3) & (data < 0.8)]
```

```
Out[9]: b    0.50
        c    0.75
        dtype: float64
```

```
In [10]: # fancy indexing
data[['a', 'e']]
```

```
Out[10]: a    0.25
         e    1.25
         dtype: float64
```

Among these, slicing may be the source of the most confusion. Notice that when slicing with an explicit index (i.e., `data['a':'c']`), the final index is *included* in the slice, while when slicing with an implicit index (i.e., `data[0:2]`), the final index is *excluded* from the slice.

Indexers: loc, iloc, and ix

These 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 [11]: data = pd.Series(['a', 'b', 'c'], index=[1, 3, 5])
data
```

```
Out[11]: 1    a
         3    b
         5    c
         dtype: object
```

```
In [12]: # explicit index when indexing
data[1]
```

```
Out[12]: 'a'
```

```
In [13]: # implicit index when slicing
data[1:3]
```

```
Out[13]: 3    b
         5    c
         dtype: object
```

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 [14]: data.loc[1]
```

```
Out[14]: 'a'
```

```
In [15]: data.loc[1:3]
```

```
Out[15]: 1    a  
        3    b  
        dtype: object
```

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

```
In [16]: data.iloc[1]
```

```
Out[16]: 'b'
```

```
In [17]: data.iloc[1:3]
```

```
Out[17]: 3    b  
        5    c  
        dtype: object
```

A third indexing attribute, `ix`, is a hybrid of the two, and for `Series` objects is equivalent to standard `[]`-based indexing. The purpose of the `ix` indexer will become more apparent in the context of `DataFrame` objects, which we will discuss in a moment.

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 [14]: 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
```

Out[14]:

	area	pop
California	423967	38332521
Texas	695662	26448193
New York	141297	19651127
Florida	170312	19552860
Illinois	149995	12882135

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

```
In [15]: data['area']
```

```
Out[15]: California    423967
Texas                695662
New York             141297
Florida              170312
Illinois             149995
Name: area, dtype: int64
```

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

```
In [16]: data.area
```

```
Out[16]: California    423967
Texas                695662
New York             141297
Florida              170312
Illinois             149995
Name: area, dtype: int64
```

This attribute-style column access actually accesses the exact same object as the dictionary-style access:

```
In [17]: data.area is data['area']
```

```
Out[17]: True
```

Though this is a useful shorthand, keep in mind that it does not work for all cases! For example, if the column names are not strings, or if the column names conflict with methods of the

`DataFrame` , this attribute-style access is not possible. For example, the `DataFrame` has a `pop()` method, so `data.pop` will point to this rather than the "pop" column:

```
In [22]: data.pop is data['pop']
```

```
Out[22]: False
```

In particular, you should avoid the temptation to try column assignment via attribute (i.e., use `data['pop'] = z` rather than `data.pop = z`).

Like with the `Series` objects discussed earlier, this dictionary-style syntax can also be used to modify the object, in this case adding a new column:

```
In [23]: data['density'] = data['pop'] / data['area']
data
```

```
Out[23]:
```

	area	pop	density
California	423967	38332521	90.413926
Texas	695662	26448193	38.018740
New York	141297	19651127	139.076746
Florida	170312	19552860	114.806121
Illinois	149995	12882135	85.883763

This shows a preview of the straightforward syntax of element-by-element arithmetic between `Series` objects; we'll dig into this further in [Operating on Data in Pandas \(03.03-Operations-in-Pandas.ipynb\)](#).

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 [18]: data.values
```

```
Out[18]: array([[ 423967, 38332521],
 [ 695662, 26448193],
 [ 141297, 19651127],
 [ 170312, 19552860],
 [ 149995, 12882135]], dtype=int64)
```

With this picture in mind, many familiar array-like observations can be done on the `DataFrame` itself. For example, we can transpose the full `DataFrame` to swap rows and columns:

```
In [19]: data.T
```

```
Out[19]:
```

	California	Texas	New York	Florida	Illinois
area	423967	695662	141297	170312	149995
pop	38332521	26448193	19651127	19552860	12882135

When it comes to indexing of `DataFrame` objects, however, it is clear that the dictionary-style indexing of columns precludes our ability to simply treat it as a NumPy array. In particular, passing a single index to an array accesses a row:

```
In [20]: data.values[0]
```

```
Out[20]: array([ 423967, 38332521], dtype=int64)
```

and passing a single "index" to a `DataFrame` accesses a column:

```
In [21]: data['area']
```

```
Out[21]: California    423967
Texas                695662
New York             141297
Florida              170312
Illinois             149995
Name: area, dtype: int64
```

Thus for array-style indexing, we need another convention. Here Pandas again uses the `loc`, `iloc`, and `ix` indexers mentioned earlier. 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 [32]: data.iloc[:3, :2]
```

```
Out[32]:
```

	area	pop
California	423967	38332521
Texas	695662	26448193
New York	141297	19651127

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

```
Out[29]:
```

	area	pop
California	423967	38332521
Florida	170312	19552860
Illinois	149995	12882135

The `ix` indexer allows a hybrid of these two approaches:

```
In [30]: data.ix[:3, : 'pop']
```

```
Out[30]:
```

	area	pop
California	423967	38332521
Florida	170312	19552860
Illinois	149995	12882135

Keep in mind that for integer indices, the `ix` indexer is subject to the same potential sources of confusion as discussed for integer-indexed `Series` objects.

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

```
Out[31]:
```

	pop	density
Florida	19552860	114.806121
New York	19651127	139.076746

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

```
Out[32]:
```

	area	pop	density
California	423967	38332521	90.000000
Florida	170312	19552860	114.806121
Illinois	149995	12882135	85.883763
New York	141297	19651127	139.076746
Texas	695662	26448193	38.018740

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.

Additional indexing conventions

There are a couple extra indexing conventions that might seem at odds with the preceding discussion, but nevertheless can be very useful in practice. First, while *indexing* refers to columns, *slicing* refers to rows:

```
In [22]: data['Florida':'Illinois']
```

Out[22]:

	area	pop
Florida	170312	19552860
Illinois	149995	12882135

Such slices can also refer to rows by number rather than by index:

```
In [34]: data[1:3]
```

Out[34]:

	area	pop	density
Florida	170312	19552860	114.806121
Illinois	149995	12882135	85.883763

Similarly, direct masking operations are also interpreted row-wise rather than column-wise:

```
In [35]: data[data.density > 100]
```

Out[35]:

	area	pop	density
Florida	170312	19552860	114.806121
New York	141297	19651127	139.076746

These two conventions are syntactically similar to those on a NumPy array, and while these may not precisely fit the mold of the Pandas conventions, they are nevertheless quite useful in practice.