

pandas

- pandas contains **data structures** and **data manipulation** tools designed to make **data cleaning and analysis fast and easy** in Python.
- pandas is often used in tandem with numerical computing tools like NumPy and SciPy, analytical libraries like statsmodels and scikit-learn, and data visualization libraries like matplotlib.
- pandas adopts significant parts of NumPy's idiomatic style of array-based computing, especially array-based functions and a preference for data processing **without for loops**.
- While pandas adopts many coding idioms from NumPy, the biggest difference is that pandas is designed for working with **tabular** or **heterogeneous** data. NumPy, by contrast, is best suited for working with homogeneous numerical array data.

We use the following import convention for pandas

```
In [ ]: import pandas as pd
```

You may also find it easier to import Series and DataFrame into the local namespace since they are so frequently used:

```
In [ ]: from pandas import Series, DataFrame
```

```
In [ ]: import numpy as np
```

Introduction to pandas Data Structures

Series

A Series is a one-dimensional array-like object containing a sequence of values (of similar types to NumPy types) and an associated array of data labels, called its index. The simplest Series is formed from only an array of data:

```
In [ ]: obj = pd.Series([4, 7, -5, 3])
obj
```

- The string representation of a Series displayed interactively shows the index on the left and the values on the right.
- Since we did not specify an index for the data, a default one consisting of the integers 0 through $N - 1$ (where N is the length of the data) is created.

```
In [ ]: obj.values
```

```
In [ ]: obj.index # like range(4)
```

Often it will be desirable to create a Series with an index identifying each data point with a label:

```
In [ ]: obj2 = pd.Series([4, 7, -5, 3], index=['d', 'b', 'a', 'c'])
obj2
```

```
In [ ]: obj2.index
```

Compared with NumPy arrays, you can use labels in the index when selecting single values or a set of values:

```
In [ ]: obj2['a']
```

```
In [ ]: obj2['d'] = 6
```

```
In [ ]: obj2[['c', 'a', 'd']]
```

Here `['c', 'a', 'd']` is interpreted as a list of indices, even though it contains strings instead of integers.

Using NumPy functions or NumPy-like operations, such as filtering with a boolean array, scalar multiplication, or applying math functions, will preserve the index-value link:

```
In [ ]: obj2[obj2 > 0]
```

```
In [ ]: obj2 * 2
```

```
In [ ]: np.exp(obj2)
```

A Series can be thought of as a fixed-length, ordered dict, as it is a mapping of index values to data values. It can be used in many contexts where you might use a dict:

```
In [ ]: 'b' in obj2
```

```
In [ ]: 'e' in obj2
```

Should you have data contained in a Python dict, you can create a Series from it by passing the dict:

```
In [ ]: sdata = {'Ohio': 35000, 'Texas': 71000, 'Oregon': 16000, 'Utah': 5000}
```

```
In [ ]: obj3 = pd.Series(sdata)
```

```
In [ ]: obj3
```

You can override the order by passing the dict keys in the order you want them to appear in the resulting Series:

```
In [ ]: states = ['California', 'Ohio', 'Oregon', 'Texas']
```

```
In [ ]: obj4 = pd.Series(sdata, index=states)
```

```
In [ ]: obj4
```

- Since no value for 'California' was found, it appears as `NaN` (not a number), which is considered in pandas to mark missing or NA values.
- Since 'Utah' was not included in states, it is excluded from the resulting object.

The `isnull` and `notnull` functions in pandas can be used to detect missing data:

```
In [ ]: pd.isnull(obj4)
```

```
In [ ]: pd.notnull(obj4)
```

Series also has these as instance methods:

```
In [ ]: obj4.isnull()
```

A useful Series feature for many applications is that it automatically aligns by index label in arithmetic operations:

```
In [ ]: obj3
```

```
In [ ]: obj4
```

```
In [ ]: obj3 + obj4
```

Data alignment features can be as a join operation of database.

Both the Series object itself and its index have a `name` attribute, which integrates with other key areas of pandas functionality:

```
In [ ]: obj4.name = 'population'
```

```
In [ ]: obj4.index.name = 'state'
```

```
In [ ]: obj4
```

A Series's index can be altered in-place by assignment:

```
In [ ]: obj
```

```
In [ ]: obj.index = ['Bob', 'Steve', 'Jeff', 'Ryan']
```

```
In [ ]: obj
```

DataFrame

- A DataFrame represents a **rectangular table of data**.
- Each of the columns can be a different value type (numeric, string, boolean, etc.).
- The DataFrame has **both a row and column index**;
 - it can be thought of as a dict of Series all sharing the same index.
 - Under the hood, the data is stored as one or more two-dimensional blocks rather than a list, dict, or some other collection of one-dimensional arrays.
 - While a DataFrame is physically two-dimensional, you can use it to represent higher dimensional data in a tabular format using **hierarchical indexing** (advanced feature of pandas).

```
In [ ]: data = {'state': ['Ohio', 'Ohio', 'Ohio', 'Nevada', 'Nevada', 'Nevada'],  
               'year': [2000, 2001, 2002, 2001, 2002, 2003],  
               'pop': [1.5, 1.7, 3.6, 2.4, 2.9, 3.2]}  
frame = pd.DataFrame(data)
```

The resulting DataFrame will have its index assigned automatically as with Series:

```
In [ ]: frame
```

For large DataFrames, the `head` method selects only the first five rows:

```
In [ ]: frame.head()
```

If you specify a sequence of columns, the DataFrame's columns will be arranged in that order:

```
In [ ]: pd.DataFrame(data, columns=['year', 'state', 'pop'])
```

If you pass a column that isn't contained in the dict, it will appear with missing values in the result:

```
In [ ]: frame2 = pd.DataFrame(data, columns=['year', 'state', 'pop', 'debt'],
                             index=['one', 'two', 'three', 'four',
                                    'five', 'six'])
frame2
```

A column in a DataFrame can be retrieved as a Series either by dict-like notation or by attribute:

```
In [ ]: frame2.columns
```

```
In [ ]: frame2['state']
```

- Attribute-like access, such as, `frame2.year` and tab completion of column names in IPython is provided as a convenience.
- `frame2[column]` works for any column name, but `frame2.column` only works when the column name is a valid Python variable name.

```
In [ ]: frame2.year
```

```
In [ ]: frame2
```

Rows can also be retrieved by position or name with the special `loc` attribute

```
In [ ]: frame2.loc['three']
```

Columns can be modified by assignment. For example, the empty 'debt' column could be assigned a scalar value or an array of values:

```
In [ ]: frame2['debt'] = 16.5
frame2
```

```
In [ ]: frame2['debt'] = np.arange(6.)
frame2
```

When you are assigning lists or arrays to a column, the value's length must match the length of the DataFrame.

If you assign a Series, its labels will be realigned exactly to the DataFrame's index, inserting missing values (i.e., NaN) in any holes:

```
In [ ]: val = pd.Series([-1.2, -1.5, -1.7], index=['two', 'four', 'five'])
val
```

```
In [ ]: frame2['debt'] = val
```

```
In [ ]: frame2
```

Assigning a column that doesn't exist will create a new column. The `del` keyword will delete columns as with a dict.

To demonstrate `del`, let us first add a new column of boolean values where the state column equals 'Ohio':

```
In [ ]: frame2['eastern'] = frame2.state == 'Ohio'
```

```
In [ ]: frame2
```

```
In [ ]: del frame2['eastern']
```

```
In [ ]: frame2.columns
```

```
In [ ]: frame2
```

The column returned from indexing a DataFrame is a view on the underlying data, not a copy.

Thus, any in-place modifications to the Series will be reflected in the DataFrame.

The column can be explicitly copied with the Series's `copy` method.

Another common form of data is a nested dict of dicts:

```
In [ ]: pop = {'Nevada': {2001: 2.4, 2002: 2.9},
               'Ohio': {2000: 1.5, 2001: 1.7, 2002: 3.6}}

In [ ]: frame3 = pd.DataFrame(pop)
frame3
```

If the nested dict is passed to the DataFrame, pandas will interpret the outer dict keys as the columns and the inner keys as the row indices.

You can transpose the DataFrame using `.T`:

```
In [ ]: frame3.T
```

The keys in the inner dicts are combined.

Index can be specified explicitly:

```
In [ ]: pd.DataFrame(pop, index=[2001, 2002, 2003])

In [ ]: frame3
```

Dicts of Series are treated in much the same way:

```
In [ ]: pdata = {'Ohio': frame3['Ohio'][:-1],
                 'Nevada': frame3['Nevada'][:2]}
pd.DataFrame(pdata)
```

For a complete list of things you can pass the DataFrame constructor, see Table 1.

Table 1: Possible data inputs to DataFrame constructor

Type	Notes
2D ndarray	A matrix of data, passing optional row and column labels
dict of arrays, lists, or tuples	Each sequence becomes a column in the DataFrame; all sequences must be the same length
NumPy structured/record array	Treated as the “dict of arrays” case
dict of Series	Each value becomes a column; indexes from each Series are unioned together to form the result’s row index if no explicit index is passed
dict of dicts	Each inner dict becomes a column; keys are unioned to form the row index as in the “dict of Series” case
List of dicts or Series	Each item becomes a row in the DataFrame; union of dict keys or Series indexes become the DataFrame’s column labels
List of lists or tuples	Treated as the “2D ndarray” case
Another DataFrame	The DataFrame’s indexes are used unless different ones are passed
NumPy MaskedArray	Like the “2D ndarray” case except masked values become NA/missing in the DataFrame result

If a DataFrame’s `index` and `columns` have their name attributes set, these will also be displayed:

```
In [ ]: frame3
```

```
In [ ]: frame3.index.name = 'year'; frame3.columns.name = 'state'  
frame3
```

As with Series, the values attribute returns the data contained in the DataFrame as a two-dimensional ndarray:

```
In [ ]: frame3.values
```

If the DataFrame's columns are different dtypes, the dtype of the values array will be chosen to accommodate all of the columns:

```
In [ ]: frame2.values
```

Index Objects

pandas's Index objects are responsible for holding the axis labels and other metadata (like the axis name or names). Any array or other sequence of labels you use when constructing a Series or DataFrame is internally converted to an Index:

```
In [ ]: obj = pd.Series(range(3), index=['a', 'b', 'c'])  
obj
```

```
In [ ]: index = obj.index
```

```
In [ ]: index
```

```
In [ ]: index[1:]
```

Index objects are immutable and thus can't be modified by the user:

```
In [ ]: index[1] = 'd' # TypeError
```

Immutability makes it safer to share Index objects among data structures:

```
In [ ]: labels = pd.Index(np.arange(3))  
labels
```

```
In [ ]: obj2 = pd.Series([1.5, -2.5, 0], index=labels)
```

```
In [ ]: obj2
```

```
In [ ]: obj2.index is labels
```

In addition to being array-like, an Index also behaves like a fixed-size set:

```
In [ ]: frame3
```

```
In [ ]: frame3.columns
```

```
In [ ]: 'Ohio' in frame3.columns
```

```
In [ ]: 2003 in frame3.index
```

Unlike Python sets, a pandas Index can contain duplicate labels:

```
In [ ]: dup_labels = pd.Index(['foo', 'foo', 'bar', 'bar'])
```

```
In [ ]: dup_labels
```

Selections with duplicate labels will select all occurrences of that label.

Each Index has a number of methods and properties for set logic, which answer other common questions about the data it contains. Some useful ones are summarized in Table 2.

Table 2: Some Index methods and properties.

Method	Description
<code>append</code>	Concatenate with additional Index objects, producing a new Index
<code>difference</code>	Compute set difference as an Index
<code>intersection</code>	Compute set intersection
<code>union</code>	Compute set union
<code>isin</code>	Compute boolean array indicating whether each value is contained in the passed collection
<code>delete</code>	Compute new Index with element at index <code>i</code> deleted
<code>drop</code>	Compute new Index by deleting passed values
<code>insert</code>	Compute new Index by inserting element at index <code>i</code>
<code>is_monotonic</code>	Returns <code>True</code> if each element is greater than or equal to the previous element
<code>is_unique</code>	Returns <code>True</code> if the Index has no duplicate values
<code>unique</code>	Compute the array of unique values in the Index

Essential Functionality

This section will walk you through the fundamental mechanics of interacting with the data contained in a Series or DataFrame.

Reindexing

An important method on pandas objects is `reindex`, which means to create a new object with the data conformed to a new index. Consider an example:

```
In [ ]: obj = pd.Series([4.5, 7.2, -5.3, 3.6], index=['d', 'b', 'a', 'c'])
obj
```

Calling `reindex` on this Series rearranges the data according to the new index, introducing missing values if any index values were not already present:

```
In [ ]: obj2 = obj.reindex(['a', 'b', 'c', 'd', 'e'])
obj2
```

For ordered data like **time series**, it may be desirable to do some interpolation or filling of values when reindexing. The `method` option allows us to do this, using a method such as `ffill`, which forward-fills the values:

```
In [ ]: obj3 = pd.Series(['blue', 'purple', 'yellow'], index=[0, 2, 4])
```

```
In [ ]: obj3
```

```
In [ ]: obj3.reindex(range(6), method='ffill')
```

With `DataFrame`, `reindex` can alter either the (row) index, columns, or both. When passed only a sequence, it reindexes the **rows** in the result:

```
In [ ]: frame = pd.DataFrame(np.arange(9).reshape((3, 3)),
                             index=['a', 'c', 'd'],
                             columns=['Ohio', 'Texas', 'California'])

frame
```

```
In [ ]: frame2 = frame.reindex(['a', 'b', 'c', 'd'])

frame2
```

The columns can be `reindexed` with the `columns` keyword:

```
In [ ]: frame
```

```
In [ ]: states = ['Texas', 'Utah', 'California']
frame.reindex(columns=states)
```

See Table 3 for more about the arguments to `reindex`.

Table 3: `reindex` function arguments

Argument	Description
<code>index</code>	New sequence to use as index. Can be <code>Index</code> instance or any other sequence-like Python data structure. An <code>Index</code> will be used exactly as is without any copying.
<code>method</code>	Interpolation (fill) method; <code>'ffill'</code> fills forward, while <code>'bfill'</code> fills backward.
<code>fill_value</code>	Substitute value to use when introducing missing data by reindexing.
<code>limit</code>	When forward- or backfilling, maximum size gap (in number of elements) to fill.
<code>tolerance</code>	When forward- or backfilling, maximum size gap (in absolute numeric distance) to fill for inexact matches.
<code>level</code>	Match simple <code>Index</code> on level of <code>MultilIndex</code> ; otherwise select subset of.
<code>copy</code>	If <code>True</code> , always copy underlying data even if new index is equivalent to old index; if <code>False</code> , do not copy the data when the indexes are equivalent.

You can `reindex` more concisely by label-indexing with `loc`, and many users prefer to use it exclusively:

```
In [ ]: frame.loc[['a', 'b', 'c', 'd'], states]
```

Dropping Entries from an Axis

Dropping one or more entries from an axis is easy if you already have an index array or list without those entries. As that can require a bit of munging and set logic, the `drop` method will return a new object with the indicated value or values deleted from an axis:

```
In [ ]: obj = pd.Series(np.arange(5.), index=['a', 'b', 'c', 'd', 'e'])

obj
```

```
In [ ]: new_obj = obj.drop('c')
```

```
In [ ]: new_obj
```

```
In [ ]: obj.drop(['d', 'c'])
```

```
In [ ]: obj # NOTE the output, why? See 'inplace' after a few cells.
```

With `DataFrame`, index values can be deleted from either axis. To illustrate this, we first create an example `DataFrame`:


```
In [ ]: data = pd.DataFrame(np.arange(16).reshape((4, 4)),
                             index=['Ohio', 'Colorado', 'Utah', 'New York'],
                             columns=['one', 'two', 'three', 'four'])

data
```

Calling `drop` with a sequence of labels will drop values from the row labels (axis 0):

```
In [ ]: data.drop(['Colorado', 'Ohio'])
```

You can drop values from the columns by passing `axis=1` or `axis='columns'` :

```
In [ ]: data.drop('two', axis=1)
```

```
In [ ]: data.drop(['two', 'four'], axis='columns')
```

Many functions, like `drop`, which modify the size or shape of a Series or DataFrame, can manipulate an object `inplace` **without returning a new object**:

```
In [ ]: obj
```

```
In [ ]: obj.drop('c', inplace=True)
```

```
In [ ]: obj #Be careful with the inplace, as it destroys any data that is dropped.
```

Indexing, Selection, and Filtering

Series indexing (`obj[...]`) works analogously to NumPy array indexing, except you can use the Series's index values instead of only integers. Here are some examples of this:

```
In [ ]: obj = pd.Series(np.arange(4.), index=['a', 'b', 'c', 'd'])
obj
```

```
In [ ]: obj['b']
```

```
In [ ]: obj[1]
```

```
In [ ]: obj[2:4]
```

```
In [ ]: obj[['b', 'a', 'd']]
```

```
In [ ]: obj[[1, 3]]
```

```
In [ ]: obj[obj < 2]
```

Important: Slicing with labels behaves differently than normal Python slicing in that the endpoint is **inclusive**:

```
In [ ]: obj
```

```
In [ ]: obj['b':'c']
```

Setting using these methods modifies the corresponding section of the Series:

```
In [ ]: obj['b':'c'] = 5
obj
```

Indexing into a DataFrame is for retrieving one or more columns either with a single value or sequence:

```
In [ ]: data = pd.DataFrame(np.arange(16).reshape((4, 4)),
                             index=['Ohio', 'Colorado', 'Utah', 'New York'],
                             columns=['one', 'two', 'three', 'four'])
data
```

```
In [ ]: data['two']
```

```
In [ ]: data[['three', 'one']]
```

Indexing like this has a few special cases. First, slicing or selecting data with a boolean array:

```
In [ ]: data[:2]
```

```
In [ ]: data[data['three'] > 5]
```

```
In [ ]: data
```

The row selection syntax `data[:2]` is provided as a convenience. Passing a single element or a list to the `[]` operator selects columns.

Another use case is in indexing with a boolean DataFrame, such as one produced by a scalar comparison:

```
In [ ]: data < 5
```

```
In [ ]: data[data < 5] = 0
```

```
In [ ]: data
```

This makes DataFrame syntactically more like a two-dimensional NumPy array in this particular case.

Selection with loc and iloc

For DataFrame label-indexing on the rows, there are two special indexing operators: `loc` and `iloc`. They enable you to select a subset of the rows and columns from a DataFrame with NumPy-like notation using either axis labels (`loc`) or integers (`iloc`).

```
In [ ]: data
```

As a preliminary example, let's select a single row and multiple columns by label:

```
In [ ]: data.loc['Colorado', ['two', 'three']]
```

We'll then perform some similar selections with integers using `iloc`:

```
In [ ]: data.iloc[2, [3, 0, 1]]
```

```
In [ ]: data.iloc[2]
```

```
In [ ]: data
```

```
In [ ]: data.iloc[[1, 2], [3, 0, 1]]
```

Both indexing functions work with slices in addition to single labels or lists of labels:

```
In [ ]: data.loc[:, 'Utah', 'two']

In [ ]: data

In [ ]: data.iloc[:, :3]

In [ ]: data.iloc[:, :3][data.three > 5]
```

So there are many ways to select and rearrange the data contained in a pandas object. For DataFrame, Table 4 provides a short summary of many of them.

Table 4: Indexing options with DataFrame.

Type	Notes
df[val]	Select single column or sequence of columns from the DataFrame; special case conveniences: boolean array (filter rows), slice (slice rows), or boolean DataFrame (set values based on some criterion)
df.loc[val]	Selects single row or subset of rows from the DataFrame by label
df.loc[:, val]	Selects single column or subset of columns by label
df.loc[val1, val2]	Select both rows and columns by label
df.iloc[where]	Selects single row or subset of rows from the DataFrame by integer position
df.iloc[:, where]	Selects single column or subset of columns by integer position
df.iloc[where_i, where_j]	Select both rows and columns by integer position
df.at[label_i, label_j]	Select a single scalar value by row and column label
df.iat[i, j]	Select a single scalar value by row and column position (integers)
reindex method	Select either rows or columns by labels
get_value, set_value methods	Select single value by row and column label

Integer Indexes

Working with pandas objects indexed by integers is something that often trips up new users due to some differences with indexing semantics on built-in Python data structures like lists and tuples. For example, you might not expect the following code to generate an error:

```
In [ ]: ser = pd.Series(np.arange(3.))

In [ ]: ser

In [ ]: ser[-1]
```

In this case, pandas could “fall back” on integer indexing, but it's difficult to do this in general without introducing subtle bugs. Here we have an index containing 0, 1, 2, but inferring what the user wants (label-based indexing or position-based) is difficult:

```
In [ ]: ser = pd.Series(np.arange(3.))

In [ ]: ser
```

On the other hand, with a non-integer index, there is no potential for ambiguity:

```
In [ ]: ser2 = pd.Series(np.arange(3.), index=['a', 'b', 'c'])

In [ ]: ser2
```

```
In [ ]: ser2[-1]
```

```
In [ ]: ser[:1]
```

To keep things consistent, if you have an axis index containing integers, data selection will always be label-oriented. For more precise handling, use `loc` (for labels) or `iloc` (for integers):

```
In [ ]: ser.loc[:2]  #Note: '2' is inclusive since 'loc' is for labels
```

Important Note: that contrary to usual python slices, **both** the start and the stop are included for `loc`

```
In [ ]: ser.iloc[:1]
```

Arithmetic and Data Alignment

An important pandas feature for some applications is the behavior of arithmetic between objects with different indexes. When you are adding together objects, if any index pairs are not the same, the respective index in the result will be the union of the index pairs. For users with database experience, this is similar to an automatic outer join on the index labels. Let's look at an example:

```
In [ ]: s1 = pd.Series([7.3, -2.5, 3.4, 1.5], index=['a', 'c', 'd', 'e'])
s1
```

```
In [ ]: s2 = pd.Series([-2.1, 3.6, -1.5, 4, 3.1],
                      index=['a', 'c', 'e', 'f', 'g'])
s2
```

Adding these together yields:

```
In [ ]: s1 + s2
```

The internal data alignment introduces missing values in the label locations that don't overlap. Missing values will then propagate in further arithmetic computations.

In the case of `DataFrame`, alignment is performed on both the rows and the columns:

```
In [ ]: df1 = pd.DataFrame(np.arange(9.).reshape((3, 3)), columns=list('bcd'),
                          index=['Ohio', 'Texas', 'Colorado'])
df1
```

```
In [ ]: df2 = pd.DataFrame(np.arange(12.).reshape((4, 3)), columns=list('bde'),
                          index=['Utah', 'Ohio', 'Texas', 'Oregon'])
df2
```

Adding these together returns a `DataFrame` whose index and columns are the unions of the ones in each `DataFrame`:

```
In [ ]: df1 + df2
```

Since the `'c'` and `'e'` columns are not found in both `DataFrame` objects, they appear as all missing in the result. The same holds for the rows whose labels are not common to both objects.

If you add `DataFrame` objects with no column or row labels in common, the result will contain all nulls:

```
In [ ]: df1 = pd.DataFrame({'A': [1, 2]})
df1
```

```
In [ ]: df2 = pd.DataFrame({'B': [3, 4]})
df2
```

```
In [ ]: df1 - df2
```

Arithmetic methods with fill values

In arithmetic operations between differently indexed objects, you might want to fill with a special value, like 0, when an axis label is found in one object but not the other:

```
In [ ]: df1 = pd.DataFrame(np.arange(12.).reshape((3, 4)),
                           columns=list('abcd'))
df1
```

```
In [ ]: df2 = pd.DataFrame(np.arange(20.).reshape((4, 5)),
                           columns=list('abcde'))
df2
```

```
In [ ]: df2.loc[1, 'b'] = np.nan
```

```
In [ ]: df2
```

Adding these together results in NA values in the locations that don't overlap:

```
In [ ]: df1 + df2
```

Using the `add` method on `df1`, you can pass `df2` and an argument to `fill_value`:

```
In [ ]: df1.add(df2, fill_value=0)
```

See Table 5 for a listing of Series and DataFrame methods for arithmetic.

Table 5: Flexible arithmetic methods.

Method	Description
<code>add</code> , <code>radd</code>	Methods for addition (+)
<code>sub</code> , <code>rsub</code>	Methods for subtraction (-)
<code>div</code> , <code>rdiv</code>	Methods for division (/)
<code>floordiv</code> , <code>rfloordiv</code>	Methods for floor division (//)
<code>mul</code> , <code>rmul</code>	Methods for multiplication (*)
<code>pow</code> , <code>rpow</code>	Methods for exponentiation (**)

Each of them has a counterpart, starting with the letter `r`, that has arguments flipped. So these two statements are equivalent:

```
In [ ]: 1 / df1
```

```
In [ ]: df1.rdiv(1)
```

Relatedly, when reindexing a Series or DataFrame, you can also specify a different `fill_value`:

```
In [ ]: df1.reindex(columns=df2.columns, fill_value=777)
```

Operations between DataFrame and Series

As with NumPy arrays of different dimensions, arithmetic between DataFrame and Series is also defined. First, as a motivating example, consider the difference between a two-dimensional array and one of its rows:

```
In [ ]: arr = np.arange(12.).reshape((3, 4))
arr
```

```
In [ ]: arr[0]
```

```
In [ ]: arr - arr[0]
```

When we subtract `arr[0]` from `arr`, the subtraction is performed once **for each row**. This is referred to as ***broadcasting***.

Operations between a DataFrame and a Series are similar:

```
In [ ]: frame = pd.DataFrame(np.arange(12.).reshape((4, 3)),
                             columns=list('bde'),
                             index=['Utah', 'Ohio', 'Texas', 'Oregon'])
frame
```

```
In [ ]: series = frame.iloc[0]
series
```

```
In [ ]: frame - series # Note: The subtraction is performed once for each row.
```

If an index value is not found in either the DataFrame's columns or the Series's index, the objects will be reindexed to form the union:

```
In [ ]: series2 = pd.Series(range(3), index=['b', 'e', 'f'])
series2
```

```
In [ ]: frame
```

```
In [ ]: frame + series2
```

If you want to instead broadcast over the columns, matching on the rows, you have to use one of the arithmetic methods. For example:

```
In [ ]: series3 = frame['d']
```

```
In [ ]: series3
```

```
In [ ]: frame
```

```
In [ ]: frame.sub(series3, axis='index')
```

The axis number that you pass is the axis to match on. In this case we mean to match on the DataFrame's row index (`axis='index'` or `axis=0`) and broadcast across.

Function Application and Mapping

NumPy ufuncs (element-wise array methods) also work with pandas objects:

```
In [ ]: frame = pd.DataFrame(np.random.randn(4, 3), columns=list('bde'),
                             index=['Utah', 'Ohio', 'Texas', 'Oregon'])
frame
```

```
In [ ]: np.abs(frame)
```

Another frequent operation is applying a function on one-dimensional arrays to each column or row. DataFrame's `apply` method does exactly this:

```
In [ ]: f = lambda x: x.max() - x.min()
        frame.apply(f)
```

Here the function f , which computes the difference between the maximum and minimum of a Series, is invoked once on each **column** in `frame`. The result is a Series having the columns of `frame` as its index.

If you pass `axis='columns'` to `apply`, the function will be invoked once per **row** instead:

```
In [ ]: frame.apply(f, axis='columns')
```

Many of the most common array statistics (like `sum` and `mean`) are DataFrame methods, so using `apply` is not necessary.

The function passed to `apply` need not return a scalar value; it can also return a Series with multiple values:

```
In [ ]: def f(x):
        return pd.Series([x.min(), x.max()], index=['min', 'max'])
        frame.apply(f)
```

Element-wise Python functions can be used, too. Suppose you wanted to compute a formatted string from each floating-point value in `frame`. You can do this with `applymap`:

```
In [ ]: format = lambda x: '%.2f' %x
        frame.applymap(format)
```

The reason for the name `applymap` is that Series has a `map` method for applying an element-wise function:

```
In [ ]: frame['e'].map(format)
```

Sorting and Ranking

Sorting a dataset by some criterion is another important built-in operation. To sort lexicographically by row or column index, use the `sort_index` method, which returns a new, sorted object:

```
In [ ]: obj = pd.Series(range(4), index=['d', 'a', 'b', 'c'])
```

```
In [ ]: obj
```

```
In [ ]: obj.sort_index()
```

With a DataFrame, you can **sort by index** on either axis:

```
In [ ]: frame = pd.DataFrame(np.arange(8).reshape((2, 4)),
                             index=['three', 'one'],
                             columns=['d', 'a', 'b', 'c'])
```

```
In [ ]: frame
```

```
In [ ]: frame.sort_index()
```

```
In [ ]: frame.sort_index(axis=1)
```

The data is sorted in ascending order by default, but can be sorted in descending order, too:

```
In [ ]: frame.sort_index(axis=1, ascending=False)
```

To sort a Series **by its values**, use its `sort_values` method:

```
In [ ]: obj = pd.Series([4, 7, -3, 2])
```

```
In [ ]: obj
```

```
In [ ]: obj.sort_values()
```

Any missing values are sorted to the end of the Series by default:

```
In [ ]: obj = pd.Series([4, np.nan, 7, np.nan, -3, 2])
```

```
In [ ]: obj
```

```
In [ ]: obj.sort_values()
```

When sorting a DataFrame, you can use the data in one or more columns as the sort keys. To do so, pass one or more column names to the `by` option of `sort_values`:

```
In [ ]: frame = pd.DataFrame({'b': [4, 7, -3, 2], 'a': [0, 1, 0, 1]})
frame
```

```
In [ ]: frame.sort_values(by='b')
```

To sort by multiple columns, pass a list of names:

```
In [ ]: frame.sort_values(by=['a', 'b'])
```

Ranking assigns ranks from one through the number of valid data points in an array. The `rank` methods for Series and DataFrame are the place to look; by default `rank` breaks ties by assigning each group the mean rank:

```
In [ ]: obj = pd.Series([7, -5, 7, 4, 2, 0, 4])
```

```
In [ ]: obj
```

```
In [ ]: obj.rank()
```

Ranks can also be assigned according to the order in which they're observed in the data.

Here, instead of using the average rank 6.5 for the entries 0 and 2, they instead have been set to 6 and 7 because label 0 precedes label 2 in the data.

```
In [ ]: obj.rank(method='first')
```

You can rank in descending order, too:

```
In [ ]: # Assign tie values the maximum rank in the group
obj.rank(ascending=False, method='max')
```


See Table 6 for a list of tie-breaking methods available.

Table 6: Tie-breaking methods with rank.

Method	Description
'average'	Default: assign the average rank to each entry in the equal group
'min'	Use the minimum rank for the whole group
'max'	Use the maximum rank for the whole group
'first'	Assign ranks in the order the values appear in the data
'dense'	Like method='min', but ranks always increase by 1 in between groups rather than the number of equal elements in a group

DataFrame can compute ranks over the rows or the columns:

```
In [ ]: frame = pd.DataFrame({'b': [4.3, 7, -3, 2], 'a': [0, 1, 0, 1],
                             'c': [-2, 5, 8, -2.5]})
frame

In [ ]: frame.rank(axis='columns')
```

Axis Indexes with Duplicate Labels

While many pandas functions (like `reindex`) require that the labels be unique, it's not mandatory. Let's consider a small Series with duplicate indices:

```
In [ ]: obj = pd.Series(range(5), index=['a', 'a', 'b', 'b', 'c'])
obj
```

The index's `is_unique` property can tell you whether its labels are unique or not:

```
In [ ]: obj.index.is_unique
```

Data selection is one of the main things that behaves differently with duplicates. Indexing a label with multiple entries returns a Series, while single entries return a scalar value:

```
In [ ]: obj['a']

In [ ]: obj['c']
```

This can make your code more complicated, as the output type from indexing can vary based on whether a label is repeated or not.

The same logic extends to indexing rows in a DataFrame:

```
In [ ]: df = pd.DataFrame(np.random.randn(4, 3), index=['a', 'a', 'b', 'b'])
df

In [ ]: df.loc['b']
```

Summarizing and Computing Descriptive Statistics

- pandas objects are equipped with a set of common mathematical and statistical methods.
- Most of these fall into the category of reductions or summary statistics, methods that extract a single value (like the `sum` or `mean`) from a Series or a Series of values from the rows or columns of a DataFrame.
- Compared with the similar methods found on NumPy arrays, they have built-in handling for **missing data**.

Consider a small DataFrame:

```
In [ ]: df = pd.DataFrame([[1.4, np.nan], [7.1, -4.5],
                           [np.nan, np.nan], [0.75, -1.3]],
                           index=['a', 'b', 'c', 'd'],
                           columns=['one', 'two'])

df
```

Calling DataFrame's `sum` method returns a Series containing column sums:

```
In [ ]: df.sum()
```

Passing `axis='columns'` or `axis=1` sums across the columns instead:

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

NA values are excluded unless the entire slice (row or column in this case) is NA. This can be disabled with the `skipna` option:

```
In [ ]: df.mean(axis='columns', skipna=False)
```

See Table 7 for a list of common options for each reduction method.

Table 7: Options for reduction methods.

Method	Description
<code>axis</code>	Axis to reduce over; 0 for DataFrame's rows and 1 for columns
<code>skipna</code>	Exclude missing values; True by default
<code>level</code>	Reduce grouped by level if the axis is hierarchically indexed (MultiIndex)

```
In [ ]: df
```

Some methods, like `idxmin` and `idxmax`, return indirect statistics like the index value where the minimum or maximum values are attained:

```
In [ ]: df.idxmax()
```

Other methods are accumulations:

```
In [ ]: df.cumsum()
```

Another type of method is neither a reduction nor an accumulation. `describe` is one such example, producing multiple summary statistics in one shot:

```
In [ ]: df.describe()
```

```
In [ ]: obj = pd.Series(['a', 'a', 'b', 'c'] * 4)
```

```
In [ ]: obj

In [ ]: # For object data (e.g. strings or timestamps), outputs are: count, unique, top, freq
# Here: The 'top' is the most common value. The 'freq' is the most common value's frequency.
obj.describe()
```

See Table 8 for a full list of summary statistics and related methods.

Table 8: Descriptive and summary statistics.

Method	Description
count	Number of non-NA values
describe	Compute set of summary statistics for Series or each DataFrame column
min, max	Compute minimum and maximum values
argmin, argmax	Compute index locations (integers) at which minimum or maximum value obtained, respectively
idxmin, idxmax	Compute index labels at which minimum or maximum value obtained, respectively
quantile	Compute sample quantile ranging from 0 to 1
sum	Sum of values
mean	Mean of values
median	Arithmetic median (50% quantile) of values
mad	Mean absolute deviation from mean value
prod	Product of all values
var	Sample variance of values
std	Sample standard deviation of values
skew	Sample skewness (third moment) of values
kurt	Sample kurtosis (fourth moment) of values
cumsum	Cumulative sum of values
cummin, cummax	Cumulative minimum or maximum of values, respectively
cumprod	Cumulative product of values
diff	Compute first arithmetic difference (useful for time series)
pct_change	Compute percent changes

Correlation and Covariance

Some summary statistics, like correlation and covariance, are computed from pairs of arguments.

```
In [ ]: # Load iris dataset [you may review the code we used in chapter 1]
import pandas as pd
iris=pd.read_csv("iris.arff")

iriscp=iris.copy()
myreplacementlist= {"Iris-setosa":0, "Iris-versicolor":1,"Iris-virginica":2}
iriscp.replace({'class ': myreplacementlist}, inplace=True)

iriscp

In [ ]: iriscp.columns
```

The `corr` method of Series computes the correlation of the overlapping, non-NA, aligned-by-index values in two Series. Relatedly, `cov` computes the covariance:

```
In [ ]: # Correlation shows whether and how strongly pairs of variables are related.

        iriscp['sepal_length'].corr(iriscp['class ']) # correlation between 'sepal_length' and 'c
        lass '
```

```
In [ ]: # covariance is a measure of the joint variability of two random variables

        iriscp['sepal_length'].cov(iriscp['class ']) # covariance:'sepal_length' versus 'class '
```

Since `sepal_length` is a valid Python attribute (but not `'class '` because of the space at the end), we can also select these columns using more concise syntax:

```
In [ ]: iriscp.sepal_length.corr(iriscp['class '])
```

`DataFrame`'s `corr` and `cov` methods, on the other hand, return a full correlation or covariance matrix as a `DataFrame`, respectively

```
In [ ]: iriscp.corr()
```

```
In [ ]: iriscp.cov()
```

Using `DataFrame`'s `corrwith` method, you can compute pairwise correlations between a `DataFrame`'s columns or rows with another Series or `DataFrame`. Passing a Series returns a Series with the correlation value computed for each column:

```
In [ ]: iriscp.corrwith(iriscp.sepal_length)
```

Unique Values, Value Counts, and Membership

Another class of related methods extracts information about the values contained in a one-dimensional Series. To illustrate these, consider this example:

```
In [ ]: obj = pd.Series(['c', 'a', 'd', 'a', 'a', 'b', 'b', 'c', 'c'])
```

The first function is `unique`, which gives you an array of the unique values in a Series:

```
In [ ]: uniques = obj.unique()
```

```
In [ ]: uniques
```

The unique values are not necessarily returned in sorted order, but could be sorted after the fact if needed (`uniques.sort()`). Relatedly, `value_counts` computes a Series containing value frequencies:

```
In [ ]: obj.value_counts()
```

The Series is sorted by value in descending order as a convenience. `value_counts` is also available as a top-level pandas method that can be used with any array or sequence:

```
In [ ]: pd.value_counts(obj.values, sort=False)
```

`isin` performs a vectorized set membership check and can be useful in filtering a dataset down to a subset of values in a Series or column in a `DataFrame`:

```
In [ ]: obj
```

```
In [ ]: mask = obj.isin(['b', 'c'])
```

```
In [ ]: mask
```

```
In [ ]: obj[mask]
```

Related to `isin` is the `Index.get_indexer` method, which gives you an index array from an array of possibly non-distinct values into another array of distinct values:

```
In [ ]: to_match = pd.Series(['c', 'a', 'b', 'b', 'c', 'a'])
unique_vals = pd.Series(['c', 'b', 'a'])
pd.Index(unique_vals).get_indexer(to_match)
```

See Table 9 for a reference on these methods.

Table 9: Unique, value counts, and set membership methods.

Method	Description
<code>isin</code>	Compute boolean array indicating whether each Series value is contained in the passed sequence of values
<code>get_indexer</code>	Compute integer indices for each value in an array into another array of distinct values; helpful for data alignment and join-type operations
<code>unique</code>	Compute array of unique values in a Series, returned in the order observed
<code>value_counts</code>	Return a Series containing unique values as its index and frequencies as its values, ordered count in descending order

In some cases, you may want to compute a histogram on multiple related columns in a DataFrame. Here's an example:

```
In [ ]: data = pd.DataFrame({'Qu1': [1, 3, 4, 3, 4],
                             'Qu2': [2, 3, 1, 2, 3],
                             'Qu3': [1, 5, 2, 4, 4]})
data
```

Passing `pandas.value_counts` to this DataFrame's `apply` function gives:

```
In [ ]: result = data.apply(pd.value_counts).fillna(0)
result
```

Here, the row labels in the result are the distinct values occurring in all of the columns. The values are the respective counts of these values in each column.

References:

- [1] Python for Data Analysis Data Wrangling with Pandas, NumPy, and IPython by Wes McKinney, 2nd Edn, O'Reilly 2017.