

### 2.1.17 SAS formats

The top-level function `read_sas()` can read (but not write) SAS *xport* (.XPT) and (since v0.18.0) *SAS7BDAT* (.sas7bdat) format files.

SAS files only contain two value types: ASCII text and floating point values (usually 8 bytes but sometimes truncated). For xport files, there is no automatic type conversion to integers, dates, or categoricals. For SAS7BDAT files, the format codes may allow date variables to be automatically converted to dates. By default the whole file is read and returned as a `DataFrame`.

Specify a `chunksize` or use `iterator=True` to obtain reader objects (`XportReader` or `SAS7BDATReader`) for incrementally reading the file. The reader objects also have attributes that contain additional information about the file and its variables.

Read a SAS7BDAT file:

```
df = pd.read_sas('sas_data.sas7bdat')
```

Obtain an iterator and read an XPORT file 100,000 lines at a time:

```
def do_something(chunk):
    pass

rdr = pd.read_sas('sas_xport.xpt', chunk=100000)
for chunk in rdr:
    do_something(chunk)
```

The [specification](#) for the xport file format is available from the SAS web site.

No official documentation is available for the SAS7BDAT format.

### 2.1.18 SPSS formats

New in version 0.25.0.

The top-level function `read_spss()` can read (but not write) SPSS *sav* (.sav) and *zsav* (.zsav) format files.

SPSS files contain column names. By default the whole file is read, categorical columns are converted into `pd.Categorical`, and a `DataFrame` with all columns is returned.

Specify the `usecols` parameter to obtain a subset of columns. Specify `convert_categoricals=False` to avoid converting categorical columns into `pd.Categorical`.

Read an SPSS file:

```
df = pd.read_spss('spss_data.sav')
```

Extract a subset of columns contained in `usecols` from an SPSS file and avoid converting categorical columns into `pd.Categorical`:

```
df = pd.read_spss('spss_data.sav', usecols=['foo', 'bar'],
                  convert_categoricals=False)
```

More information about the *sav* and *zsav* file format is available [here](#).

### 2.1.19 Other file formats

pandas itself only supports IO with a limited set of file formats that map cleanly to its tabular data model. For reading and writing other file formats into and from pandas, we recommend these packages from the broader community.

#### netCDF

`xarray` provides data structures inspired by the pandas `DataFrame` for working with multi-dimensional datasets, with a focus on the netCDF file format and easy conversion to and from pandas.

### 2.1.20 Performance considerations

This is an informal comparison of various IO methods, using pandas 0.24.2. Timings are machine dependent and small differences should be ignored.

```
In [1]: sz = 1000000
In [2]: df = pd.DataFrame({'A': np.random.randn(sz), 'B': [1] * sz})

In [3]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000000 entries, 0 to 999999
Data columns (total 2 columns):
A      1000000 non-null float64
B      1000000 non-null int64
dtypes: float64(1), int64(1)
memory usage: 15.3 MB
```

Given the next test set:

```
import numpy as np

import os

sz = 1000000
df = pd.DataFrame({'A': np.random.randn(sz), 'B': [1] * sz})

sz = 1000000
np.random.seed(42)
df = pd.DataFrame({'A': np.random.randn(sz), 'B': [1] * sz})

def test_sql_write(df):
    if os.path.exists('test.sql'):
        os.remove('test.sql')
    sql_db = sqlite3.connect('test.sql')
    df.to_sql(name='test_table', con=sql_db)
    sql_db.close()

def test_sql_read():
    sql_db = sqlite3.connect('test.sql')
    pd.read_sql_query("select * from test_table", sql_db)
    sql_db.close()

def test_hdf_fixed_write(df):
    df.to_hdf('test_fixed.hdf', 'test', mode='w')
```

(continues on next page)

(continued from previous page)

```

def test_hdf_fixed_read():
    pd.read_hdf('test_fixed.hdf', 'test')

def test_hdf_fixed_write_compress(df):
    df.to_hdf('test_fixed_compress.hdf', 'test', mode='w', complib='blosc')

def test_hdf_fixed_read_compress():
    pd.read_hdf('test_fixed_compress.hdf', 'test')

def test_hdf_table_write(df):
    df.to_hdf('test_table.hdf', 'test', mode='w', format='table')

def test_hdf_table_read():
    pd.read_hdf('test_table.hdf', 'test')

def test_hdf_table_write_compress(df):
    df.to_hdf('test_table_compress.hdf', 'test', mode='w',
              complib='blosc', format='table')

def test_hdf_table_read_compress():
    pd.read_hdf('test_table_compress.hdf', 'test')

def test_csv_write(df):
    df.to_csv('test.csv', mode='w')

def test_csv_read():
    pd.read_csv('test.csv', index_col=0)

def test_feather_write(df):
    df.to_feather('test.feather')

def test_feather_read():
    pd.read_feather('test.feather')

def test_pickle_write(df):
    df.to_pickle('test.pkl')

def test_pickle_read():
    pd.read_pickle('test.pkl')

def test_pickle_write_compress(df):
    df.to_pickle('test.pkl.compress', compression='xz')

def test_pickle_read_compress():
    pd.read_pickle('test.pkl.compress', compression='xz')

def test_parquet_write(df):
    df.to_parquet('test.parquet')

def test_parquet_read():
    pd.read_parquet('test.parquet')

```

When writing, the top-three functions in terms of speed are `test_feather_write`, `test_hdf_fixed_write` and `test_hdf_fixed_write_compress`.

```
In [4]: %timeit test_sql_write(df)
```

(continues on next page)

(continued from previous page)

```

3.29 s ± 43.2 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [5]: %timeit test_hdf_fixed_write(df)
19.4 ms ± 560 µs per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [6]: %timeit test_hdf_fixed_write_compress(df)
19.6 ms ± 308 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [7]: %timeit test_hdf_table_write(df)
449 ms ± 5.61 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [8]: %timeit test_hdf_table_write_compress(df)
448 ms ± 11.9 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [9]: %timeit test_csv_write(df)
3.66 s ± 26.2 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [10]: %timeit test_feather_write(df)
9.75 ms ± 117 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)

In [11]: %timeit test_pickle_write(df)
30.1 ms ± 229 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [12]: %timeit test_pickle_write_compress(df)
4.29 s ± 15.9 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [13]: %timeit test_parquet_write(df)
67.6 ms ± 706 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

```

When reading, the top three are `test_feather_read`, `test_pickle_read` and `test_hdf_fixed_read`.

```

In [14]: %timeit test_sql_read()
1.77 s ± 17.7 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [15]: %timeit test_hdf_fixed_read()
19.4 ms ± 436 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [16]: %timeit test_hdf_fixed_read_compress()
19.5 ms ± 222 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [17]: %timeit test_hdf_table_read()
38.6 ms ± 857 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [18]: %timeit test_hdf_table_read_compress()
38.8 ms ± 1.49 ms per loop (mean ± std. dev. of 7 runs, 10 loops each)

In [19]: %timeit test_csv_read()
452 ms ± 9.04 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

In [20]: %timeit test_feather_read()
12.4 ms ± 99.7 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)

In [21]: %timeit test_pickle_read()
18.4 ms ± 191 µs per loop (mean ± std. dev. of 7 runs, 100 loops each)

In [22]: %timeit test_pickle_read_compress()
915 ms ± 7.48 ms per loop (mean ± std. dev. of 7 runs, 1 loop each)

```

(continues on next page)

(continued from previous page)

```
In [23]: %timeit test_parquet_read()
24.4 ms ± 146 µs per loop (mean ± std. dev. of 7 runs, 10 loops each)
```

For this test case `test.pkl.compress`, `test.parquet` and `test.feather` took the least space on disk. Space on disk (in bytes)

```
29519500 Oct 10 06:45 test.csv
16000248 Oct 10 06:45 test.feather
8281983 Oct 10 06:49 test.parquet
16000857 Oct 10 06:47 test.pkl
7552144 Oct 10 06:48 test.pkl.compress
34816000 Oct 10 06:42 test.sql
24009288 Oct 10 06:43 test_fixed.hdf
24009288 Oct 10 06:43 test_fixed_compress.hdf
24458940 Oct 10 06:44 test_table.hdf
24458940 Oct 10 06:44 test_table_compress.hdf
```

## 2.2 Indexing and selecting data

The axis labeling information in pandas objects serves many purposes:

- Identifies data (i.e. provides *metadata*) using known indicators, important for analysis, visualization, and interactive console display.
- Enables automatic and explicit data alignment.
- Allows intuitive getting and setting of subsets of the data set.

In this section, we will focus on the final point: namely, how to slice, dice, and generally get and set subsets of pandas objects. The primary focus will be on Series and DataFrame as they have received more development attention in this area.

**Note:** The Python and NumPy indexing operators `[]` and attribute operator `.` provide quick and easy access to pandas data structures across a wide range of use cases. This makes interactive work intuitive, as there's little new to learn if you already know how to deal with Python dictionaries and NumPy arrays. However, since the type of the data to be accessed isn't known in advance, directly using standard operators has some optimization limits. For production code, we recommended that you take advantage of the optimized pandas data access methods exposed in this chapter.

**Warning:** Whether a copy or a reference is returned for a setting operation, may depend on the context. This is sometimes called `chained assignment` and should be avoided. See [Returning a View versus Copy](#).

See the [MultiIndex / Advanced Indexing](#) for MultiIndex and more advanced indexing documentation.

See the [cookbook](#) for some advanced strategies.

## 2.2.1 Different choices for indexing

Object selection has had a number of user-requested additions in order to support more explicit location based indexing. Pandas now supports three types of multi-axis indexing.

- `.loc` is primarily label based, but may also be used with a boolean array. `.loc` will raise `KeyError` when the items are not found. Allowed inputs are:
  - A single label, e.g. 5 or 'a' (Note that 5 is interpreted as a *label* of the index. This use is **not** an integer position along the index.).
  - A list or array of labels ['a', 'b', 'c'].
  - A slice object with labels 'a': 'f' (Note that contrary to usual python slices, **both** the start and the stop are included, when present in the index! See [Slicing with labels](#) and [Endpoints are inclusive](#).)
  - A boolean array (any NA values will be treated as `False`).
  - A callable function with one argument (the calling Series or DataFrame) and that returns valid output for indexing (one of the above).

See more at [Selection by Label](#).

- `.iloc` is primarily integer position based (from 0 to `length-1` of the axis), but may also be used with a boolean array. `.iloc` will raise `IndexError` if a requested indexer is out-of-bounds, except *slice* indexers which allow out-of-bounds indexing. (this conforms with Python/NumPy *slice* semantics). Allowed inputs are:
  - An integer e.g. 5.
  - A list or array of integers [4, 3, 0].
  - A slice object with ints 1:7.
  - A boolean array (any NA values will be treated as `False`).
  - A callable function with one argument (the calling Series or DataFrame) and that returns valid output for indexing (one of the above).

See more at [Selection by Position](#), [Advanced Indexing](#) and [Advanced Hierarchical](#).

- `.loc`, `.iloc`, and also `[]` indexing can accept a callable as indexer. See more at [Selection By Callable](#).

Getting values from an object with multi-axes selection uses the following notation (using `.loc` as an example, but the following applies to `.iloc` as well). Any of the axes accessors may be the null slice `:`. Axes left out of the specification are assumed to be `:`, e.g. `p.loc['a']` is equivalent to `p.loc['a', :, :]`.

Object Type	Indexers
Series	<code>s.loc[indexer]</code>
DataFrame	<code>df.loc[row_indexer, column_indexer]</code>

## 2.2.2 Basics

As mentioned when introducing the data structures in the [last section](#), the primary function of indexing with `[]` (a.k.a. `__getitem__` for those familiar with implementing class behavior in Python) is selecting out lower-dimensional slices. The following table shows return type values when indexing pandas objects with `[]`:

Object Type	Selection	Return Value Type
Series	<code>series[label]</code>	scalar value
DataFrame	<code>frame[colname]</code>	Series corresponding to colname

Here we construct a simple time series data set to use for illustrating the indexing functionality:

```
In [1]: dates = pd.date_range('1/1/2000', periods=8)

In [2]: df = pd.DataFrame(np.random.randn(8, 4),
    ...:                   index=dates, columns=['A', 'B', 'C', 'D'])
    ...:

In [3]: df
Out[3]:
```

	A	B	C	D
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804
2000-01-04	0.721555	-0.706771	-1.039575	0.271860
2000-01-05	-0.424972	0.567020	0.276232	-1.087401
2000-01-06	-0.673690	0.113648	-1.478427	0.524988
2000-01-07	0.404705	0.577046	-1.715002	-1.039268
2000-01-08	-0.370647	-1.157892	-1.344312	0.844885

**Note:** None of the indexing functionality is time series specific unless specifically stated.

Thus, as per above, we have the most basic indexing using []:

```
In [4]: s = df['A']

In [5]: s[dates[5]]
Out[5]: -0.6736897080883706
```

You can pass a list of columns to [] to select columns in that order. If a column is not contained in the DataFrame, an exception will be raised. Multiple columns can also be set in this manner:

```
In [6]: df
Out[6]:
```

	A	B	C	D
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804
2000-01-04	0.721555	-0.706771	-1.039575	0.271860
2000-01-05	-0.424972	0.567020	0.276232	-1.087401
2000-01-06	-0.673690	0.113648	-1.478427	0.524988
2000-01-07	0.404705	0.577046	-1.715002	-1.039268
2000-01-08	-0.370647	-1.157892	-1.344312	0.844885

```
In [7]: df[['B', 'A']] = df[['A', 'B']]

In [8]: df
Out[8]:
```

	A	B	C	D
2000-01-01	-0.282863	0.469112	-1.509059	-1.135632
2000-01-02	-0.173215	1.212112	0.119209	-1.044236
2000-01-03	-2.104569	-0.861849	-0.494929	1.071804
2000-01-04	-0.706771	0.721555	-1.039575	0.271860
2000-01-05	0.567020	-0.424972	0.276232	-1.087401
2000-01-06	0.113648	-0.673690	-1.478427	0.524988
2000-01-07	0.577046	0.404705	-1.715002	-1.039268

(continues on next page)

(continued from previous page)

```
2000-01-08 -1.157892 -0.370647 -1.344312  0.844885
```

You may find this useful for applying a transform (in-place) to a subset of the columns.

**Warning:** pandas aligns all AXES when setting Series and DataFrame from `.loc`, and `.iloc`.

This will **not** modify `df` because the column alignment is before value assignment.

```
In [9]: df[['A', 'B']]
```

```
Out[9]:
```

```

      A      B
2000-01-01 -0.282863  0.469112
2000-01-02 -0.173215  1.212112
2000-01-03 -2.104569 -0.861849
2000-01-04 -0.706771  0.721555
2000-01-05  0.567020 -0.424972
2000-01-06  0.113648 -0.673690
2000-01-07  0.577046  0.404705
2000-01-08 -1.157892 -0.370647
```

```
In [10]: df.loc[:, ['B', 'A']] = df[['A', 'B']]
```

```
In [11]: df[['A', 'B']]
```

```
Out[11]:
```

```

      A      B
2000-01-01 -0.282863  0.469112
2000-01-02 -0.173215  1.212112
2000-01-03 -2.104569 -0.861849
2000-01-04 -0.706771  0.721555
2000-01-05  0.567020 -0.424972
2000-01-06  0.113648 -0.673690
2000-01-07  0.577046  0.404705
2000-01-08 -1.157892 -0.370647
```

The correct way to swap column values is by using raw values:

```
In [12]: df.loc[:, ['B', 'A']] = df[['A', 'B']].to_numpy()
```

```
In [13]: df[['A', 'B']]
```

```
Out[13]:
```

```

      A      B
2000-01-01  0.469112 -0.282863
2000-01-02  1.212112 -0.173215
2000-01-03 -0.861849 -2.104569
2000-01-04  0.721555 -0.706771
2000-01-05 -0.424972  0.567020
2000-01-06 -0.673690  0.113648
2000-01-07  0.404705  0.577046
2000-01-08 -0.370647 -1.157892
```



### 2.2.3 Attribute access

You may access an index on a Series or column on a DataFrame directly as an attribute:

```
In [14]: sa = pd.Series([1, 2, 3], index=list('abc'))
```

```
In [15]: dfa = df.copy()
```

```
In [16]: sa.b
```

```
Out[16]: 2
```

```
In [17]: dfa.A
```

```
Out[17]:
```

```
2000-01-01    0.469112
```

```
2000-01-02    1.212112
```

```
2000-01-03   -0.861849
```

```
2000-01-04    0.721555
```

```
2000-01-05   -0.424972
```

```
2000-01-06   -0.673690
```

```
2000-01-07    0.404705
```

```
2000-01-08   -0.370647
```

```
Freq: D, Name: A, dtype: float64
```

```
In [18]: sa.a = 5
```

```
In [19]: sa
```

```
Out[19]:
```

```
a      5
```

```
b      2
```

```
c      3
```

```
dtype: int64
```

```
In [20]: dfa.A = list(range(len(dfa.index))) # ok if A already exists
```

```
In [21]: dfa
```

```
Out[21]:
```

```
          A          B          C          D
```

```
2000-01-01  0 -0.282863 -1.509059 -1.135632
```

```
2000-01-02  1 -0.173215  0.119209 -1.044236
```

```
2000-01-03  2 -2.104569 -0.494929  1.071804
```

```
2000-01-04  3 -0.706771 -1.039575  0.271860
```

```
2000-01-05  4  0.567020  0.276232 -1.087401
```

```
2000-01-06  5  0.113648 -1.478427  0.524988
```

```
2000-01-07  6  0.577046 -1.715002 -1.039268
```

```
2000-01-08  7 -1.157892 -1.344312  0.844885
```

```
In [22]: dfa['A'] = list(range(len(dfa.index))) # use this form to create a new_
↪column
```

```
In [23]: dfa
```

```
Out[23]:
```

```
          A          B          C          D
```

```
2000-01-01  0 -0.282863 -1.509059 -1.135632
```

```
2000-01-02  1 -0.173215  0.119209 -1.044236
```

```
2000-01-03  2 -2.104569 -0.494929  1.071804
```

```
2000-01-04  3 -0.706771 -1.039575  0.271860
```

```
2000-01-05  4  0.567020  0.276232 -1.087401
```

(continues on next page)

(continued from previous page)

```
2000-01-06  5  0.113648 -1.478427  0.524988
2000-01-07  6  0.577046 -1.715002 -1.039268
2000-01-08  7 -1.157892 -1.344312  0.844885
```

**Warning:**

- You can use this access only if the index element is a valid Python identifier, e.g. `s.1` is not allowed. See [here for an explanation of valid identifiers](#).
- The attribute will not be available if it conflicts with an existing method name, e.g. `s.min` is not allowed, but `s['min']` is possible.
- Similarly, the attribute will not be available if it conflicts with any of the following list: `index`, `major_axis`, `minor_axis`, `items`.
- In any of these cases, standard indexing will still work, e.g. `s['1']`, `s['min']`, and `s['index']` will access the corresponding element or column.

If you are using the IPython environment, you may also use tab-completion to see these accessible attributes.

You can also assign a dict to a row of a DataFrame:

```
In [24]: x = pd.DataFrame({'x': [1, 2, 3], 'y': [3, 4, 5]})
```

```
In [25]: x.iloc[1] = {'x': 9, 'y': 99}
```

```
In [26]: x
```

```
Out [26]:
```

```
   x  y
0  1  3
1  9 99
2  3  5
```

You can use attribute access to modify an existing element of a Series or column of a DataFrame, but be careful; if you try to use attribute access to create a new column, it creates a new attribute rather than a new column. In 0.21.0 and later, this will raise a UserWarning:

```
In [1]: df = pd.DataFrame({'one': [1., 2., 3.]})
```

```
In [2]: df.two = [4, 5, 6]
```

```
UserWarning: Pandas doesn't allow Series to be assigned into nonexistent columns -
↳ see https://pandas.pydata.org/pandas-docs/stable/indexing.html#attribute_access
```

```
In [3]: df
```

```
Out [3]:
```

```
   one
0  1.0
1  2.0
2  3.0
```

## 2.2.4 Slicing ranges

The most robust and consistent way of slicing ranges along arbitrary axes is described in the [Selection by Position](#) section detailing the `.iloc` method. For now, we explain the semantics of slicing using the `[]` operator.

With Series, the syntax works exactly as with an ndarray, returning a slice of the values and the corresponding labels:

```
In [27]: s[:5]
Out[27]:
2000-01-01    0.469112
2000-01-02    1.212112
2000-01-03   -0.861849
2000-01-04    0.721555
2000-01-05   -0.424972
Freq: D, Name: A, dtype: float64

In [28]: s[::2]
Out[28]:
2000-01-01    0.469112
2000-01-03   -0.861849
2000-01-05   -0.424972
2000-01-07    0.404705
Freq: 2D, Name: A, dtype: float64

In [29]: s[::-1]
Out[29]:
2000-01-08   -0.370647
2000-01-07    0.404705
2000-01-06   -0.673690
2000-01-05   -0.424972
2000-01-04    0.721555
2000-01-03   -0.861849
2000-01-02    1.212112
2000-01-01    0.469112
Freq: -1D, Name: A, dtype: float64
```

Note that setting works as well:

```
In [30]: s2 = s.copy()

In [31]: s2[:5] = 0

In [32]: s2
Out[32]:
2000-01-01    0.000000
2000-01-02    0.000000
2000-01-03    0.000000
2000-01-04    0.000000
2000-01-05    0.000000
2000-01-06   -0.673690
2000-01-07    0.404705
2000-01-08   -0.370647
Freq: D, Name: A, dtype: float64
```

With DataFrame, slicing inside of `[]` **slices the rows**. This is provided largely as a convenience since it is such a common operation.

```
In [33]: df[:3]
Out[33]:
```

	A	B	C	D
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804

```
In [34]: df[::-1]
Out[34]:
```

	A	B	C	D
2000-01-08	-0.370647	-1.157892	-1.344312	0.844885
2000-01-07	0.404705	0.577046	-1.715002	-1.039268
2000-01-06	-0.673690	0.113648	-1.478427	0.524988
2000-01-05	-0.424972	0.567020	0.276232	-1.087401
2000-01-04	0.721555	-0.706771	-1.039575	0.271860
2000-01-03	-0.861849	-2.104569	-0.494929	1.071804
2000-01-02	1.212112	-0.173215	0.119209	-1.044236
2000-01-01	0.469112	-0.282863	-1.509059	-1.135632

## 2.2.5 Selection by label

**Warning:** Whether a copy or a reference is returned for a setting operation, may depend on the context. This is sometimes called `chained assignment` and should be avoided. See [Returning a View versus Copy](#).

### Warning:

`.loc` is strict when you present slicers that are not compatible (or convertible) with the index type. For example using integers in a `DatetimeIndex`. These will raise a `TypeError`.

```
In [35]: df1 = pd.DataFrame(np.random.randn(5, 4),
.....:                      columns=list('ABCD'),
.....:                      index=pd.date_range('20130101', periods=5))
.....:
```

```
In [36]: df1
```

```
Out[36]:
```

	A	B	C	D
2013-01-01	1.075770	-0.109050	1.643563	-1.469388
2013-01-02	0.357021	-0.674600	-1.776904	-0.968914
2013-01-03	-1.294524	0.413738	0.276662	-0.472035
2013-01-04	-0.013960	-0.362543	-0.006154	-0.923061
2013-01-05	0.895717	0.805244	-1.206412	2.565646

```
In [4]: df1.loc[2:3]
TypeError: cannot do slice indexing on <class 'pandas.tseries.index.DatetimeIndex'>
↪with these indexers [2] of <type 'int'>
```

String likes in slicing *can* be convertible to the type of the index and lead to natural slicing.

```
In [37]: df1.loc['20130102':'20130104']
Out[37]:
```

	A	B	C	D
2013-01-02	0.357021	-0.674600	-1.776904	-0.968914
2013-01-03	-1.294524	0.413738	0.276662	-0.472035
2013-01-04	-0.013960	-0.362543	-0.006154	-0.923061

**Warning:** Starting in 0.21.0, pandas will show a `FutureWarning` if indexing with a list with missing labels. In the future this will raise a `KeyError`. See *list-like Using loc with missing keys in a list is Deprecated*.

pandas provides a suite of methods in order to have **purely label based indexing**. This is a strict inclusion based protocol. Every label asked for must be in the index, or a `KeyError` will be raised. When slicing, both the start bound **AND** the stop bound are *included*, if present in the index. Integers are valid labels, but they refer to the label **and not the position**.

The `.loc` attribute is the primary access method. The following are valid inputs:

- A single label, e.g. 5 or 'a' (Note that 5 is interpreted as a *label* of the index. This use is **not** an integer position along the index.).
- A list or array of labels ['a', 'b', 'c'].
- A slice object with labels 'a': 'f' (Note that contrary to usual python slices, **both** the start and the stop are included, when present in the index! See *Slicing with labels*).
- A boolean array.
- A callable, see *Selection By Callable*.

```
In [38]: s1 = pd.Series(np.random.randn(6), index=list('abcdef'))
```

```
In [39]: s1
```

```
Out[39]:
```

```
a    1.431256
b    1.340309
c   -1.170299
d   -0.226169
e    0.410835
f    0.813850
dtype: float64
```

```
In [40]: s1.loc['c:']
```

```
Out[40]:
```

```
c   -1.170299
d   -0.226169
e    0.410835
f    0.813850
dtype: float64
```

```
In [41]: s1.loc['b']
```

```
Out[41]: 1.3403088497993827
```

Note that setting works as well:

```
In [42]: s1.loc['c:'] = 0
```

```
In [43]: s1
```

```
Out[43]:
```

```
a    1.431256
b    1.340309
c    0.000000
d    0.000000
```

(continues on next page)

(continued from previous page)

```
e    0.000000
f    0.000000
dtype: float64
```

With a DataFrame:

```
In [44]: df1 = pd.DataFrame(np.random.randn(6, 4),
.....:                      index=list('abcdef'),
.....:                      columns=list('ABCD'))
.....:
```

```
In [45]: df1
```

```
Out[45]:
```

	A	B	C	D
a	0.132003	-0.827317	-0.076467	-1.187678
b	1.130127	-1.436737	-1.413681	1.607920
c	1.024180	0.569605	0.875906	-2.211372
d	0.974466	-2.006747	-0.410001	-0.078638
e	0.545952	-1.219217	-1.226825	0.769804
f	-1.281247	-0.727707	-0.121306	-0.097883

```
In [46]: df1.loc[['a', 'b', 'd'], :]
```

```
Out[46]:
```

	A	B	C	D
a	0.132003	-0.827317	-0.076467	-1.187678
b	1.130127	-1.436737	-1.413681	1.607920
d	0.974466	-2.006747	-0.410001	-0.078638

Accessing via label slices:

```
In [47]: df1.loc['d':, 'A':'C']
```

```
Out[47]:
```

	A	B	C
d	0.974466	-2.006747	-0.410001
e	0.545952	-1.219217	-1.226825
f	-1.281247	-0.727707	-0.121306

For getting a cross section using a label (equivalent to `df.xs('a')`):

```
In [48]: df1.loc['a']
```

```
Out[48]:
```

A	0.132003
B	-0.827317
C	-0.076467
D	-1.187678

Name: a, dtype: float64

For getting values with a boolean array:

```
In [49]: df1.loc['a'] > 0
```

```
Out[49]:
```

A	True
B	False
C	False
D	False

Name: a, dtype: bool

(continues on next page)

(continued from previous page)

```
In [50]: df1.loc[:, df1.loc['a'] > 0]
Out[50]:
```

	A
a	0.132003
b	1.130127
c	1.024180
d	0.974466
e	0.545952
f	-1.281247

NA values in a boolean array propagate as False:

Changed in version 1.0.2: `mask = pd.array([True, False, True, False, pd.NA, False], dtype="boolean")` `mask df1[mask]`

For getting a value explicitly:

```
# this is also equivalent to ``df1.at['a', 'A']``
In [51]: df1.loc['a', 'A']
Out[51]: 0.13200317033032932
```

## Slicing with labels

When using `.loc` with slices, if both the start and the stop labels are present in the index, then elements *located* between the two (including them) are returned:

```
In [52]: s = pd.Series(list('abcde'), index=[0, 3, 2, 5, 4])

In [53]: s.loc[3:5]
Out[53]:
```

3	b
2	c
5	d

dtype: object

If at least one of the two is absent, but the index is sorted, and can be compared against start and stop labels, then slicing will still work as expected, by selecting labels which *rank* between the two:

```
In [54]: s.sort_index()
Out[54]:
```

0	a
2	c
3	b
4	e
5	d

dtype: object

```
In [55]: s.sort_index().loc[1:6]
Out[55]:
```

2	c
3	b
4	e
5	d

dtype: object

However, if at least one of the two is absent *and* the index is not sorted, an error will be raised (since doing otherwise would be computationally expensive, as well as potentially ambiguous for mixed type indexes). For instance, in the above example, `s.loc[1:6]` would raise `KeyError`.

For the rationale behind this behavior, see *Endpoints are inclusive*.

## 2.2.6 Selection by position

**Warning:** Whether a copy or a reference is returned for a setting operation, may depend on the context. This is sometimes called `chained assignment` and should be avoided. See *Returning a View versus Copy*.

Pandas provides a suite of methods in order to get **purely integer based indexing**. The semantics follow closely Python and NumPy slicing. These are 0-based indexing. When slicing, the start bound is *included*, while the upper bound is *excluded*. Trying to use a non-integer, even a **valid** label will raise an `IndexError`.

The `.iloc` attribute is the primary access method. The following are valid inputs:

- An integer e.g. 5.
- A list or array of integers [4, 3, 0].
- A slice object with ints 1:7.
- A boolean array.
- A callable, see *Selection By Callable*.

```
In [56]: s1 = pd.Series(np.random.randn(5), index=list(range(0, 10, 2)))

In [57]: s1
Out[57]:
0    0.695775
2    0.341734
4    0.959726
6   -1.110336
8   -0.619976
dtype: float64

In [58]: s1.iloc[:3]
Out[58]:
0    0.695775
2    0.341734
4    0.959726
dtype: float64

In [59]: s1.iloc[3]
Out[59]: -1.110336102891167
```

Note that setting works as well:

```
In [60]: s1.iloc[:3] = 0

In [61]: s1
Out[61]:
0    0.000000
2    0.000000
4    0.000000
```

(continues on next page)



(continued from previous page)

```
6    -1.110336
8    -0.619976
dtype: float64
```

With a DataFrame:

```
In [62]: df1 = pd.DataFrame(np.random.randn(6, 4),
.....:                      index=list(range(0, 12, 2)),
.....:                      columns=list(range(0, 8, 2)))
.....:
```

```
In [63]: df1
```

```
Out [63]:
```

	0	2	4	6
0	0.149748	-0.732339	0.687738	0.176444
2	0.403310	-0.154951	0.301624	-2.179861
4	-1.369849	-0.954208	1.462696	-1.743161
6	-0.826591	-0.345352	1.314232	0.690579
8	0.995761	2.396780	0.014871	3.357427
10	-0.317441	-1.236269	0.896171	-0.487602

Select via integer slicing:

```
In [64]: df1.iloc[:3]
```

```
Out [64]:
```

	0	2	4	6
0	0.149748	-0.732339	0.687738	0.176444
2	0.403310	-0.154951	0.301624	-2.179861
4	-1.369849	-0.954208	1.462696	-1.743161

```
In [65]: df1.iloc[1:5, 2:4]
```

```
Out [65]:
```

	4	6
2	0.301624	-2.179861
4	1.462696	-1.743161
6	1.314232	0.690579
8	0.014871	3.357427

Select via integer list:

```
In [66]: df1.iloc[[1, 3, 5], [1, 3]]
```

```
Out [66]:
```

	2	6
2	-0.154951	-2.179861
6	-0.345352	0.690579
10	-1.236269	-0.487602

```
In [67]: df1.iloc[1:3, :]
```

```
Out [67]:
```

	0	2	4	6
2	0.403310	-0.154951	0.301624	-2.179861
4	-1.369849	-0.954208	1.462696	-1.743161

```
In [68]: df1.iloc[:, 1:3]
```

```
Out [68]:
```

	2	4
--	---	---

(continues on next page)

(continued from previous page)

```
0 -0.732339 0.687738
2 -0.154951 0.301624
4 -0.954208 1.462696
6 -0.345352 1.314232
8 2.396780 0.014871
10 -1.236269 0.896171
```

```
# this is also equivalent to ``df1.iat[1,1]``
In [69]: df1.iloc[1, 1]
Out[69]: -0.1549507744249032
```

For getting a cross section using an integer position (equiv to `df.xs(1)`):

```
In [70]: df1.iloc[1]
Out[70]:
0    0.403310
2   -0.154951
4    0.301624
6   -2.179861
Name: 2, dtype: float64
```

Out of range slice indexes are handled gracefully just as in Python/Numpy.

```
# these are allowed in python/numpy.
In [71]: x = list('abcdef')

In [72]: x
Out[72]: ['a', 'b', 'c', 'd', 'e', 'f']

In [73]: x[4:10]
Out[73]: ['e', 'f']

In [74]: x[8:10]
Out[74]: []

In [75]: s = pd.Series(x)

In [76]: s
Out[76]:
0    a
1    b
2    c
3    d
4    e
5    f
dtype: object

In [77]: s.iloc[4:10]
Out[77]:
4    e
5    f
dtype: object

In [78]: s.iloc[8:10]
Out[78]: Series([], dtype: object)
```

Note that using slices that go out of bounds can result in an empty axis (e.g. an empty DataFrame being returned).

```
In [79]: df1 = pd.DataFrame(np.random.randn(5, 2), columns=list('AB'))
```

```
In [80]: df1
```

```
Out[80]:
```

	A	B
0	-0.082240	-2.182937
1	0.380396	0.084844
2	0.432390	1.519970
3	-0.493662	0.600178
4	0.274230	0.132885

```
In [81]: df1.iloc[:, 2:3]
```

```
Out[81]:
```

Empty DataFrame  
Columns: []  
Index: [0, 1, 2, 3, 4]

```
In [82]: df1.iloc[:, 1:3]
```

```
Out[82]:
```

	B
0	-2.182937
1	0.084844
2	1.519970
3	0.600178
4	0.132885

```
In [83]: df1.iloc[4:6]
```

```
Out[83]:
```

	A	B
4	0.27423	0.132885

A single indexer that is out of bounds will raise an `IndexError`. A list of indexers where any element is out of bounds will raise an `IndexError`.

```
>>> df1.iloc[[4, 5, 6]]
IndexError: positional indexers are out-of-bounds

>>> df1.iloc[:, 4]
IndexError: single positional indexer is out-of-bounds
```

## 2.2.7 Selection by callable

`.loc`, `.iloc`, and also `[]` indexing can accept a callable as indexer. The callable must be a function with one argument (the calling Series or DataFrame) that returns valid output for indexing.

```
In [84]: df1 = pd.DataFrame(np.random.randn(6, 4),
.....:                      index=list('abcdef'),
.....:                      columns=list('ABCD'))
.....:
```

```
In [85]: df1
```

```
Out[85]:
```

	A	B	C	D
a	-0.023688	2.410179	1.450520	0.206053
b	-0.251905	-2.213588	1.063327	1.266143

(continues on next page)

(continued from previous page)

```

c  0.299368 -0.863838  0.408204 -1.048089
d -0.025747 -0.988387  0.094055  1.262731
e  1.289997  0.082423 -0.055758  0.536580
f -0.489682  0.369374 -0.034571 -2.484478

In [86]: df1.loc[lambda df: df['A'] > 0, :]
Out[86]:
      A      B      C      D
c  0.299368 -0.863838  0.408204 -1.048089
e  1.289997  0.082423 -0.055758  0.536580

In [87]: df1.loc[:, lambda df: ['A', 'B']]
Out[87]:
      A      B
a -0.023688  2.410179
b -0.251905 -2.213588
c  0.299368 -0.863838
d -0.025747 -0.988387
e  1.289997  0.082423
f -0.489682  0.369374

In [88]: df1.iloc[:, lambda df: [0, 1]]
Out[88]:
      A      B
a -0.023688  2.410179
b -0.251905 -2.213588
c  0.299368 -0.863838
d -0.025747 -0.988387
e  1.289997  0.082423
f -0.489682  0.369374

In [89]: df1[lambda df: df.columns[0]]
Out[89]:
a   -0.023688
b   -0.251905
c    0.299368
d   -0.025747
e    1.289997
f   -0.489682
Name: A, dtype: float64

```

You can use callable indexing in Series.

```

In [90]: df1['A'].loc[lambda s: s > 0]
Out[90]:
c    0.299368
e    1.289997
Name: A, dtype: float64

```

Using these methods / indexers, you can chain data selection operations without using a temporary variable.

```

In [91]: bb = pd.read_csv('data/baseball.csv', index_col='id')

In [92]: (bb.groupby(['year', 'team']).sum()
.....:      .loc[lambda df: df['r'] > 100])
.....:
Out[92]:

```

(continues on next page)

(continued from previous page)

		stint	g	ab	r	h	X2b	X3b	hr	rbi	sb	cs	bb	so	
↪ibb	hbp	sh	sf	gidp											↪
year	team														↪
↪															
2007	CIN	6	379	745	101	203	35	2	36	125.0	10.0	1.0	105	127.0	14.
↪0	1.0	1.0	15.0	18.0											
	DET	5	301	1062	162	283	54	4	37	144.0	24.0	7.0	97	176.0	3.
↪0	10.0	4.0	8.0	28.0											
	HOU	4	311	926	109	218	47	6	14	77.0	10.0	4.0	60	212.0	3.
↪0	9.0	16.0	6.0	17.0											
	LAN	11	413	1021	153	293	61	3	36	154.0	7.0	5.0	114	141.0	8.
↪0	9.0	3.0	8.0	29.0											
	NYN	13	622	1854	240	509	101	3	61	243.0	22.0	4.0	174	310.0	24.
↪0	23.0	18.0	15.0	48.0											
	SFN	5	482	1305	198	337	67	6	40	171.0	26.0	7.0	235	188.0	51.
↪0	8.0	16.0	6.0	41.0											
	TEX	2	198	729	115	200	40	4	28	115.0	21.0	4.0	73	140.0	4.
↪0	5.0	2.0	8.0	16.0											
	TOR	4	459	1408	187	378	96	2	58	223.0	4.0	2.0	190	265.0	16.
↪0	12.0	4.0	16.0	38.0											

## 2.2.8 IX indexer is deprecated

**Warning:** Starting in 0.20.0, the `.ix` indexer is deprecated, in favor of the more strict `.iLoc` and `.loc` indexers.

`.ix` offers a lot of magic on the inference of what the user wants to do. To wit, `.ix` can decide to index *positionally* OR via *labels* depending on the data type of the index. This has caused quite a bit of user confusion over the years.

The recommended methods of indexing are:

- `.loc` if you want to *label* index.
- `.iLoc` if you want to *positionally* index.

```
In [93]: dfd = pd.DataFrame({'A': [1, 2, 3],
.....:                      'B': [4, 5, 6]},
.....:                      index=list('abc'))
.....:
```

```
In [94]: dfd
```

```
Out[94]:
```

```
   A  B
a  1  4
b  2  5
c  3  6
```

Previous behavior, where you wish to get the 0th and the 2nd elements from the index in the 'A' column.

```
In [3]: dfd.ix[[0, 2], 'A']
Out[3]:
a      1
c      3
Name: A, dtype: int64
```

Using `.loc`. Here we will select the appropriate indexes from the index, then use *label* indexing.

```
In [95]: dfd.loc[dfd.index[[0, 2]], 'A']
Out[95]:
a      1
c      3
Name: A, dtype: int64
```

This can also be expressed using `.iloc`, by explicitly getting locations on the indexers, and using *positional* indexing to select things.

```
In [96]: dfd.iloc[[0, 2], dfd.columns.get_loc('A')]
Out[96]:
a      1
c      3
Name: A, dtype: int64
```

For getting *multiple* indexers, using `.get_indexer`:

```
In [97]: dfd.iloc[[0, 2], dfd.columns.get_indexer(['A', 'B'])]
Out[97]:
   A  B
a  1  4
c  3  6
```

## 2.2.9 Indexing with list with missing labels is deprecated

**Warning:** Starting in 0.21.0, using `.loc` or `[]` with a list with one or more missing labels, is deprecated, in favor of `.reindex`.

In prior versions, using `.loc[list-of-labels]` would work as long as *at least 1* of the keys was found (otherwise it would raise a `KeyError`). This behavior is deprecated and will show a warning message pointing to this section. The recommended alternative is to use `.reindex()`.

For example.

```
In [98]: s = pd.Series([1, 2, 3])

In [99]: s
Out[99]:
0      1
1      2
2      3
dtype: int64
```

Selection with all keys found is unchanged.

```
In [100]: s.loc[[1, 2]]
Out[100]:
1      2
2      3
dtype: int64
```

Previous behavior

```
In [4]: s.loc[[1, 2, 3]]
Out[4]:
1    2.0
2    3.0
3     NaN
dtype: float64
```

#### Current behavior

```
In [4]: s.loc[[1, 2, 3]]
Passing list-likes to .loc with any non-matching elements will raise
KeyError in the future, you can use .reindex() as an alternative.

See the documentation here:
https://pandas.pydata.org/pandas-docs/stable/indexing.html#deprecate-loc-reindex-
→listlike

Out[4]:
1    2.0
2    3.0
3     NaN
dtype: float64
```

## Reindexing

The idiomatic way to achieve selecting potentially not-found elements is via `.reindex()`. See also the section on *reindexing*.

```
In [101]: s.reindex([1, 2, 3])
Out[101]:
1    2.0
2    3.0
3     NaN
dtype: float64
```

Alternatively, if you want to select only *valid* keys, the following is idiomatic and efficient; it is guaranteed to preserve the dtype of the selection.

```
In [102]: labels = [1, 2, 3]

In [103]: s.loc[s.index.intersection(labels)]
Out[103]:
1    2
2    3
dtype: int64
```

Having a duplicated index will raise for a `.reindex()`:

```
In [104]: s = pd.Series(np.arange(4), index=['a', 'a', 'b', 'c'])

In [105]: labels = ['c', 'd']
```

```
In [17]: s.reindex(labels)
ValueError: cannot reindex from a duplicate axis
```

Generally, you can intersect the desired labels with the current axis, and then reindex.

```
In [106]: s.loc[s.index.intersection(labels)].reindex(labels)
Out[106]:
c      3.0
d      NaN
dtype: float64
```

However, this would *still* raise if your resulting index is duplicated.

```
In [41]: labels = ['a', 'd']

In [42]: s.loc[s.index.intersection(labels)].reindex(labels)
ValueError: cannot reindex from a duplicate axis
```

## 2.2.10 Selecting random samples

A random selection of rows or columns from a Series or DataFrame with the `sample()` method. The method will sample rows by default, and accepts a specific number of rows/columns to return, or a fraction of rows.

```
In [107]: s = pd.Series([0, 1, 2, 3, 4, 5])

# When no arguments are passed, returns 1 row.
In [108]: s.sample()
Out[108]:
4      4
dtype: int64

# One may specify either a number of rows:
In [109]: s.sample(n=3)
Out[109]:
0      0
4      4
1      1
dtype: int64

# Or a fraction of the rows:
In [110]: s.sample(frac=0.5)
Out[110]:
5      5
3      3
1      1
dtype: int64
```

By default, `sample` will return each row at most once, but one can also sample with replacement using the `replace` option:

```
In [111]: s = pd.Series([0, 1, 2, 3, 4, 5])

# Without replacement (default):
In [112]: s.sample(n=6, replace=False)
Out[112]:
0      0
1      1
5      5
3      3
2      2
```

(continues on next page)



(continued from previous page)

```

4      4
dtype: int64

# With replacement:
In [113]: s.sample(n=6, replace=True)
Out[113]:
0      0
4      4
3      3
2      2
4      4
4      4
dtype: int64

```

By default, each row has an equal probability of being selected, but if you want rows to have different probabilities, you can pass the `sample` function sampling weights as `weights`. These weights can be a list, a NumPy array, or a Series, but they must be of the same length as the object you are sampling. Missing values will be treated as a weight of zero, and inf values are not allowed. If weights do not sum to 1, they will be re-normalized by dividing all weights by the sum of the weights. For example:

```

In [114]: s = pd.Series([0, 1, 2, 3, 4, 5])

In [115]: example_weights = [0, 0, 0.2, 0.2, 0.2, 0.4]

In [116]: s.sample(n=3, weights=example_weights)
Out[116]:
5      5
4      4
3      3
dtype: int64

# Weights will be re-normalized automatically
In [117]: example_weights2 = [0.5, 0, 0, 0, 0, 0]

In [118]: s.sample(n=1, weights=example_weights2)
Out[118]:
0      0
dtype: int64

```

When applied to a DataFrame, you can use a column of the DataFrame as sampling weights (provided you are sampling rows and not columns) by simply passing the name of the column as a string.

```

In [119]: df2 = pd.DataFrame({'col1': [9, 8, 7, 6],
.....:                       'weight_column': [0.5, 0.4, 0.1, 0]})
.....:

In [120]: df2.sample(n=3, weights='weight_column')
Out[120]:
   col1  weight_column
1      8             0.4
0      9             0.5
2      7             0.1

```

`sample` also allows users to sample columns instead of rows using the `axis` argument.

```

In [121]: df3 = pd.DataFrame({'col1': [1, 2, 3], 'col2': [2, 3, 4]})

```

(continues on next page)