Q Search the docs ...

10 minutes to pandas

Intro to data structures

Essential basic functionality

IO tools (text, CSV, HDF5, ...)

Indexing and selecting data

MultiIndex / advanced indexing

Merge, join, concatenate and compare

Reshaping and pivot tables

Working with text data

Working with missing data

Duplicate Labels

Categorical data

Nullable integer data type

Nullable Boolean data type

Chart Visualization

Table Visualization

Computational tools

Group by: split-apply-combine

Windowing Operations

Time series / date functionality

<u>Time deltas</u>

Options and settings

Enhancing performance

Scaling to large datasets

Sparse data structures

Frequently Asked Questions (FAQ)

Cookbook

10 minutes to pandas ¶

This is a short introduction to pandas, geared mainly for new users. You can see more complex recipes in the <u>Cookbook</u>.

Customarily, we import as follows:

```
In [1]: import numpy as np
In [2]: import pandas as pd
```

Object creation

See the <u>Data Structure Intro section</u>.

Creating a **Series** by passing a list of values, letting pandas create a default integer index:

```
In [3]: s = pd.Series([1, 3, 5, np.nan, 6, 8])
In [4]: s
Out[4]:
0    1.0
1    3.0
2    5.0
3    NaN
4    6.0
5    8.0
dtype: float64
```

Creating a **DataFrame** by passing a NumPy array, with a datetime index and labeled columns:

Creating a DataFrame by passing a dict of objects that can be converted to series-like.

```
In [9]: df2 = pd.DataFrame(
   • • • •
            {
                 "A": 1.0,
   . . . :
                 "B": pd.Timestamp("20130102"),
   . . . :
                 "C": pd.Series(1, index=list(range(4)), dtype="float32"),
   . . . :
                 "D": np.array([3] * 4, dtype="int32"),
   . . . :
                 "E": pd.Categorical(["test", "train", "test", "train"]),
   . . . :
                "F": "foo",
   . . . :
   ...:
   ...: )
   . . . :
In [10]: df2
Out[10]:
                     C D
                                Ε
0 1.0 2013-01-02 1.0 3
                             test
1 1.0 2013-01-02 1.0 3 train
                                   foo
2 1.0 2013-01-02 1.0 3
                                   foo
                            test
3 1.0 2013-01-02 1.0 3 train
```

The columns of the resulting **DataFrame** have different dtypes.

```
In [11]: df2.dtypes
Out[11]:
A      float64
B      datetime64[ns]
C      float32
D         int32
E         category
F         object
dtype: object
```

If you're using IPython, tab completion for column names (as well as public attributes) is automatically enabled. Here's a subset of the attributes that will be completed:

As you can see, the columns A, B, C, and D are automatically tab completed. E and F are there as well; the rest of the attributes have been truncated for brevity.

Viewing data

See the **Basics section**.

Here is how to view the top and bottom rows of the frame:

```
In [13]: df.head()
Out[13]:

A B C D

2013-01-01 0.469112 -0.282863 -1.509059 -1.135632
2013-01-02 1.212112 -0.173215 0.119209 -1.044236
2013-01-03 -0.861849 -2.104569 -0.494929 1.071804
2013-01-04 0.721555 -0.706771 -1.039575 0.271860
2013-01-05 -0.424972 0.567020 0.276232 -1.087401

In [14]: df.tail(3)
Out[14]:

A B C D

2013-01-04 0.721555 -0.706771 -1.039575 0.271860
2013-01-05 -0.424972 0.567020 0.276232 -1.087401
2013-01-06 -0.673690 0.113648 -1.478427 0.524988
```

Display the index, columns:

<u>DataFrame.to numpy()</u> gives a NumPy representation of the underlying data. Note that this can be an expensive operation when your <u>DataFrame</u> has columns with different data types, which comes down to a fundamental difference between pandas and NumPy: **NumPy arrays have one dtype for the entire array, while pandas DataFrames have one dtype per column**. When you call <u>DataFrame.to numpy()</u>, pandas will find the NumPy dtype that can hold *all* of the dtypes in the DataFrame. This may end up being object, which requires casting every value to a Python object.

For df, our <u>DataFrame</u> of all floating-point values, <u>DataFrame.to numpy()</u> is fast and doesn't require copying data.

For df2, the <u>DataFrame</u> with multiple dtypes, <u>DataFrame.to_numpy()</u> is relatively expensive.

```
In [18]: df2.to_numpy()
Out[18]:
array([[1.0, Timestamp('2013-01-02 00:00:00'), 1.0, 3, 'test', 'foo'],
        [1.0, Timestamp('2013-01-02 00:00:00'), 1.0, 3, 'train', 'foo'],
        [1.0, Timestamp('2013-01-02 00:00:00'), 1.0, 3, 'test', 'foo'],
        [1.0, Timestamp('2013-01-02 00:00:00'), 1.0, 3, 'train', 'foo']],
        dtype=object)
```

1 Note

<u>DataFrame.to_numpy()</u> does not include the index or column labels in the output.

describe() shows a quick statistic summary of your data:

Transposing your data:

```
In [20]: df.T
Out[20]:
    2013-01-01   2013-01-02   2013-01-03   2013-01-04   2013-01-05   2013-01-06
A    0.469112   1.212112   -0.861849   0.721555   -0.424972   -0.673690
B    -0.282863   -0.173215   -2.104569   -0.706771   0.567020   0.113648
C    -1.509059   0.119209   -0.494929   -1.039575   0.276232   -1.478427
D    -1.135632   -1.044236   1.071804   0.271860   -1.087401   0.524988
```

Sorting by an axis:

Sorting by values:

```
In [22]: df.sort_values(by="B")
Out[22]:

A B C D

2013-01-03 -0.861849 -2.104569 -0.494929 1.071804
2013-01-04 0.721555 -0.706771 -1.039575 0.271860
2013-01-01 0.469112 -0.282863 -1.509059 -1.135632
2013-01-02 1.212112 -0.173215 0.119209 -1.044236
2013-01-06 -0.673690 0.113648 -1.478427 0.524988
2013-01-05 -0.424972 0.567020 0.276232 -1.087401
```

Selection



While standard Python / NumPy expressions for selecting and setting are intuitive and come in handy for interactive work, for production code, we recommend the optimized pandas data access methods, .at, .iat, .loc and .iloc.

See the indexing documentation <u>Indexing and Selecting Data</u> and <u>MultiIndex / Advanced Indexing</u>.

Getting

Selecting a single column, which yields a **Series**, equivalent to df.A:

```
In [23]: df["A"]
Out[23]:
2013-01-01     0.469112
2013-01-02     1.212112
2013-01-03     -0.861849
2013-01-04     0.721555
2013-01-05     -0.424972
2013-01-06     -0.673690
Freq: D, Name: A, dtype: float64
```

Selecting via [], which slices the rows.

```
In [24]: df[0:3]
Out[24]:

A B C D

2013-01-01 0.469112 -0.282863 -1.509059 -1.135632
2013-01-02 1.212112 -0.173215 0.119209 -1.044236
2013-01-03 -0.861849 -2.104569 -0.494929 1.071804

In [25]: df["20130102":"20130104"]
Out[25]:

A B C D

2013-01-02 1.212112 -0.173215 0.119209 -1.044236
2013-01-03 -0.861849 -2.104569 -0.494929 1.071804

2013-01-04 0.721555 -0.706771 -1.039575 0.271860
```

Selection by label

See more in Selection by Label.

For getting a cross section using a label:

```
In [26]: df.loc[dates[0]]
Out[26]:
A    0.469112
B    -0.282863
C    -1.509059
D    -1.135632
Name: 2013-01-01 00:00:00, dtype: float64
```

Selecting on a multi-axis by label:

```
In [27]: df.loc[:, ["A", "B"]]
Out[27]:

A B

2013-01-01 0.469112 -0.282863

2013-01-02 1.212112 -0.173215

2013-01-03 -0.861849 -2.104569

2013-01-04 0.721555 -0.706771

2013-01-05 -0.424972 0.567020

2013-01-06 -0.673690 0.113648
```

Showing label slicing, both endpoints are included:

Reduction in the dimensions of the returned object:

```
In [29]: df.loc["20130102", ["A", "B"]]
Out[29]:
A    1.212112
B    -0.173215
Name: 2013-01-02 00:00:00, dtype: float64
```

For getting a scalar value:

```
In [30]: df.loc[dates[0], "A"]
Out[30]: 0.4691122999071863
```

For getting fast access to a scalar (equivalent to the prior method):

```
In [31]: df.at[dates[0], "A"]
Out[31]: 0.4691122999071863
```

Selection by position

See more in Selection by Position.

Select via the position of the passed integers:

```
In [32]: df.iloc[3]
Out[32]:
A   0.721555
B   -0.706771
C   -1.039575
D   0.271860
Name: 2013-01-04 00:00:00, dtype: float64
```

By integer slices, acting similar to NumPy/Python:

By lists of integer position locations, similar to the NumPy/Python style:

For slicing rows explicitly:

For slicing columns explicitly:

For getting a value explicitly:

```
In [37]: df.iloc[1, 1]
Out[37]: -0.17321464905330858
```

For getting fast access to a scalar (equivalent to the prior method):

```
In [38]: df.iat[1, 1]
Out[38]: -0.17321464905330858
```

Boolean indexing

Using a single column's values to select data.

Selecting values from a DataFrame where a boolean condition is met.

```
In [40]: df[df > 0]
Out[40]:
                                  C
                                           D
2013-01-01 0.469112
                       NaN
                                 NaN
                                          NaN
2013-01-02 1.212112
                       NaN 0.119209
                                          NaN
2013-01-03
          NaN
                       NaN
                                 NaN 1.071804
2013-01-04 0.721555
                       NaN
                                 NaN 0.271860
2013-01-05 NaN 0.567020 0.276232
                                         NaN
2013-01-06
              NaN 0.113648
                                NaN 0.524988
```

Using the <u>isin()</u> method for filtering:

```
In [41]: df2 = df.copy()
In [42]: df2["E"] = ["one", "one", "two", "three", "four", "three"]
In [43]: df2
Out[43]:
                          В
2013-01-01   0.469112   -0.282863   -1.509059   -1.135632
2013-01-02 1.212112 -0.173215 0.119209 -1.044236
2013-01-03 -0.861849 -2.104569 -0.494929 1.071804
                                                two
2013-01-05 -0.424972 0.567020 0.276232 -1.087401
2013-01-06 -0.673690 0.113648 -1.478427 0.524988 three
In [44]: df2[df2["E"].isin(["two", "four"])]
Out[44]:
                          В
2013-01-03 -0.861849 -2.104569 -0.494929 1.071804
                                               two
2013-01-05 -0.424972 0.567020 0.276232 -1.087401 four
```

Setting

Setting a new column automatically aligns the data by the indexes.

Setting values by label:

```
In [48]: df.at[dates[0], "A"] = 0
```

Setting values by position:

```
In [49]: df.iat[0, 1] = 0
```

Setting by assigning with a NumPy array:

```
In [50]: df.loc[:, "D"] = np.array([5] * len(df))
```

The result of the prior setting operations.

```
In [51]: df
Out[51]:

A B C D F

2013-01-01 0.000000 0.000000 -1.509059 5 NaN

2013-01-02 1.212112 -0.173215 0.119209 5 1.0

2013-01-03 -0.861849 -2.104569 -0.494929 5 2.0

2013-01-04 0.721555 -0.706771 -1.039575 5 3.0

2013-01-05 -0.424972 0.567020 0.276232 5 4.0

2013-01-06 -0.673690 0.113648 -1.478427 5 5.0
```

A where operation with setting.

Missing data

pandas primarily uses the value np.nan to represent missing data. It is by default not included in computations. See the <u>Missing Data section</u>.

Reindexing allows you to change/add/delete the index on a specified axis. This returns a copy of the data.

To drop any rows that have missing data.

Filling missing data.

To get the boolean mask where values are nan.

Operations

See the Basic section on Binary Ops.

Stats

Operations in general exclude missing data.

Performing a descriptive statistic:

```
In [61]: df.mean()
Out[61]:
A   -0.004474
B   -0.383981
C   -0.687758
D   5.000000
F   3.000000
dtype: float64
```

Same operation on the other axis:

```
In [62]: df.mean(1)
Out[62]:
2013-01-01    0.872735
2013-01-02    1.431621
2013-01-03    0.707731
2013-01-04    1.395042
2013-01-05    1.883656
2013-01-06    1.592306
Freq: D, dtype: float64
```

Operating with objects that have different dimensionality and need alignment. In addition, pandas automatically broadcasts along the specified dimension.

```
In [63]: s = pd.Series([1, 3, 5, np.nan, 6, 8], index=dates).shift(2)
In [64]: s
Out[64]:
2013-01-01
             NaN
2013-01-02
            NaN
2013-01-03
            1.0
2013-01-04
            3.0
2013-01-05
            5.0
2013-01-06
Freq: D, dtype: float64
In [65]: df.sub(s, axis="index")
Out[65]:
A B C D F
2013-01-01 NaN NaN NaN NaN NaN
2013-01-02 NaN NaN NaN NaN NaN
2013-01-03 -1.861849 -3.104569 -1.494929 4.0 1.0
2013-01-04 -2.278445 -3.706771 -4.039575 2.0 0.0
2013-01-05 -5.424972 -4.432980 -4.723768 0.0 -1.0
2013-01-06 NaN NaN NaN NaN NaN
```

Apply

Applying functions to the data:

```
In [66]: df.apply(np.cumsum)
Out[66]:
                      В
                             C
                                D
                                     F
              Α
2013-01-01 0.000000 0.000000 -1.509059 5
                                    NaN
2013-01-02 1.212112 -0.173215 -1.389850 10
                                    1.0
3.0
2013-01-04 1.071818 -2.984555 -2.924354 20 6.0
2013-01-06 -0.026844 -2.303886 -4.126549 30 15.0
In [67]: df.apply(lambda x: x.max() - x.min())
Out[67]:
   2.073961
Α
   2.671590
В
C
   1.785291
D
   0.000000
   4.000000
F
dtype: float64
```

Histogramming

See more at <u>Histogramming and Discretization</u>.

String Methods

Series is equipped with a set of string processing methods in the str attribute that make it easy to operate on each element of the array, as in the code snippet below. Note that pattern-matching in str generally uses <u>regular expressions</u> by default (and in some cases always uses them). See more at <u>Vectorized String Methods</u>.

```
In [71]: s = pd.Series(["A", "B", "C", "Aaba", "Baca", np.nan, "CABA", "dog", "cat"])
In [72]: s.str.lower()
Out[72]:
1
       b
2
3
    aaba
4
    baca
5
     NaN
6
    caba
7
     dog
     cat
dtype: object
```

Merge

Concat

pandas provides various facilities for easily combining together Series and DataFrame objects with various kinds of set logic for the indexes and relational algebra functionality in the case of join / merge-type operations.

See the <u>Merging section</u>.

Concatenating pandas objects together with concat():

```
In [73]: df = pd.DataFrame(np.random.randn(10, 4))
In [74]: df
Out[74]:
                             2
                   1
0 -0.548702 1.467327 -1.015962 -0.483075
1 1.637550 -1.217659 -0.291519 -1.745505
2 -0.263952 0.991460 -0.919069 0.266046
3 -0.709661 1.669052 1.037882 -1.705775
4 -0.919854 -0.042379 1.247642 -0.009920
5 0.290213 0.495767 0.362949 1.548106
6 -1.131345 -0.089329 0.337863 -0.945867
7 -0.932132 1.956030 0.017587 -0.016692
8 -0.575247 0.254161 -1.143704 0.215897
9 1.193555 -0.077118 -0.408530 -0.862495
# break it into pieces
In [75]: pieces = [df[:3], df[3:7], df[7:]]
In [76]: pd.concat(pieces)
Out[76]:
         0
                             2
                   1
0 -0.548702 1.467327 -1.015962 -0.483075
1 1.637550 -1.217659 -0.291519 -1.745505
2 -0.263952 0.991460 -0.919069 0.266046
3 -0.709661 1.669052 1.037882 -1.705775
4 -0.919854 -0.042379 1.247642 -0.009920
5 0.290213 0.495767 0.362949 1.548106
6 -1.131345 -0.089329 0.337863 -0.945867
7 -0.932132 1.956030 0.017587 -0.016692
8 -0.575247 0.254161 -1.143704 0.215897
9 1.193555 -0.077118 -0.408530 -0.862495
```

1 Note

Adding a column to a <u>DataFrame</u> is relatively fast. However, adding a row requires a copy, and may be expensive. We recommend passing a pre-built list of records to the <u>DataFrame</u> constructor instead of building a <u>DataFrame</u> by iteratively appending records to it. See <u>Appending to dataframe</u> for more.

Join

SQL style merges. See the <u>Database style joining</u> section.

```
In [77]: left = pd.DataFrame({"key": ["foo", "foo"], "lval": [1, 2]})
In [78]: right = pd.DataFrame({"key": ["foo", "foo"], "rval": [4, 5]})
In [79]: left
Out[79]:
  key lval
0 foo
          1
1 foo
In [80]: right
Out[80]:
  key rval
0 foo
          4
1 foo
In [81]: pd.merge(left, right, on="key")
Out[81]:
  key lval rval
  foo
         1
                5
1 foo
          1
2
  foo
                4
                5
3
  foo
```

Another example that can be given is:

```
In [82]: left = pd.DataFrame({"key": ["foo", "bar"], "lval": [1, 2]})
In [83]: right = pd.DataFrame({"key": ["foo", "bar"], "rval": [4, 5]})
In [84]: left
Out[84]:
  key lval
0 foo
          1
1 bar
In [85]: right
Out[85]:
  key rval
0 foo
          4
1 bar
In [86]: pd.merge(left, right, on="key")
Out[86]:
  key lval rval
0 foo
       1
1 bar
          2
```

Grouping

By "group by" we are referring to a process involving one or more of the following steps:

- **Splitting** the data into groups based on some criteria
- Applying a function to each group independently
- Combining the results into a data structure

See the **Grouping section**.

```
In [87]: df = pd.DataFrame(
   . . . . :
                   "A": ["foo", "bar", "foo", "bar", "foo", "bar", "foo", "foo"],
"B": ["one", "one", "two", "three", "two", "two", "one", "three"],
   . . . . :
   . . . . :
                   "C": np.random.randn(8),
   . . . . :
                   "D": np.random.randn(8),
   . . . . :
   . . . . :
              }
   ....: )
   . . . . :
In [88]: df
Out[88]:
                   C
     Α
            В
          one 1.346061 -1.577585
0 foo
1 bar
          one 1.511763 0.396823
          two 1.627081 -0.105381
2 foo
3 bar three -0.990582 -0.532532
4 foo
          two -0.441652 1.453749
5 bar
          two 1.211526 1.208843
6 foo
          one 0.268520 -0.080952
7 foo three 0.024580 -0.264610
```

Grouping and then applying the sum() function to the resulting groups.

Grouping by multiple columns forms a hierarchical index, and again we can apply the sum() function.

Reshaping

See the sections on <u>Hierarchical Indexing</u> and <u>Reshaping</u>.

Stack

```
In [91]: tuples = list(
            zip(
   . . . . :
   . . . . :
                     ["bar", "bar", "baz", "foo", "foo", "qux", "qux"],
   . . . . :
                     ["one", "two", "one", "two", "one", "two", "one", "two"],
   . . . . :
   . . . . :
   . . . . :
             )
   ....: )
   . . . . :
In [92]: index = pd.MultiIndex.from_tuples(tuples, names=["first", "second"])
In [93]: df = pd.DataFrame(np.random.randn(8, 2), index=index, columns=["A", "B"])
In [94]: df2 = df[:4]
In [95]: df2
Out[95]:
first second
     one
             -0.727965 -0.589346
      two
              0.339969 -0.693205
             -0.339355 0.593616
baz
     one
              0.884345 1.591431
      two
```

The stack() method "compresses" a level in the DataFrame's columns.

```
In [96]: stacked = df2.stack()
In [97]: stacked
Out[97]:
first second
bar
             A -0.727965
     one
             B -0.589346
           A 0.339969
      two
             B -0.693205
            A -0.339355
baz
     one
             B 0.593616
             A 0.884345
      two
                 1.591431
             В
dtype: float64
```

With a "stacked" DataFrame or Series (having a MultiIndex as the index), the inverse operation of stack() is unstack(), which by default unstacks the last level:

```
In [98]: stacked.unstack()
Out[98]:
first second
bar one -0.727965 -0.589346
           0.339969 -0.693205
     two
            -0.339355 0.593616
baz one
            0.884345 1.591431
     two
In [99]: stacked.unstack(1)
Out[99]:
second
             one
first
bar A -0.727965 0.339969
     B -0.589346 -0.693205
baz A -0.339355 0.884345
     B 0.593616 1.591431
In [100]: stacked.unstack(0)
Out[100]:
first
              bar
                       baz
second
    A -0.727965 -0.339355
one
      B -0.589346 0.593616
    A 0.339969 0.884345
two
      B -0.693205 1.591431
```

Pivot tables

See the section on <u>Pivot Tables</u>.

```
In [101]: df = pd.DataFrame(
                "A": ["one", "one", "two", "three"] * 3,
                "B": ["A", "B", "C"] * 4,
               "C": ["foo", "foo", "bar", "bar", "bar"] * 2,
               "D": np.random.randn(12),
                "E": np.random.randn(12),
  ····: }
  ....:)
  . . . . . :
In [102]: df
Out[102]:
            C
                   D
       A B
     one A foo -1.202872 0.047609
1
   one B foo -1.814470 -0.136473
2
    two C foo 1.018601 -0.561757
3 three A bar -0.595447 -1.623033
   one B bar 1.395433 0.029399
5
   one C bar -0.392670 -0.542108
6
   two A foo 0.007207 0.282696
7 three B foo 1.928123 -0.087302
8
  one C foo -0.055224 -1.575170
9
    one A bar 2.395985 1.771208
10
    two B bar 1.552825 0.816482
11 three C bar 0.166599 1.100230
```

We can produce pivot tables from this data very easily:

Time series

pandas has simple, powerful, and efficient functionality for performing resampling operations during frequency conversion (e.g., converting secondly data into 5-minutely data). This is extremely common in, but not limited to, financial applications. See the <u>Time Series section</u>.

Time zone representation:

```
In [107]: rng = pd.date range("3/6/2012 00:00", periods=5, freq="D")
In [108]: ts = pd.Series(np.random.randn(len(rng)), rng)
In [109]: ts
Out[109]:
            1.857704
2012-03-06
2012-03-07 -1.193545
2012-03-08
           0.677510
2012-03-09 -0.153931
2012-03-10 0.520091
Freq: D, dtype: float64
In [110]: ts_utc = ts.tz_localize("UTC")
In [111]: ts_utc
Out[111]:
2012-03-06 00:00:00+00:00
                           1.857704
2012-03-07 00:00:00+00:00 -1.193545
2012-03-08 00:00:00+00:00 0.677510
2012-03-09 00:00:00+00:00 -0.153931
2012-03-10 00:00:00+00:00
                           0.520091
Freq: D, dtype: float64
```

Converting to another time zone:

```
In [112]: ts_utc.tz_convert("US/Eastern")
Out[112]:
2012-03-05 19:00:00-05:00    1.857704
2012-03-06 19:00:00-05:00    -1.193545
2012-03-07 19:00:00-05:00    0.677510
2012-03-08 19:00:00-05:00    -0.153931
2012-03-09 19:00:00-05:00    0.520091
Freq: D, dtype: float64
```

Converting between time span representations:

```
In [113]: rng = pd.date_range("1/1/2012", periods=5, freq="M")
In [114]: ts = pd.Series(np.random.randn(len(rng)), index=rng)
In [115]: ts
Out[115]:
2012-01-31 -1.475051
           0.722570
2012-02-29
2012-03-31 -0.322646
2012-04-30 -1.601631
2012-05-31 0.778033
Freq: M, dtype: float64
In [116]: ps = ts.to_period()
In [117]: ps
Out[117]:
2012-01 -1.475051
2012-02 0.722570
2012-03 -0.322646
2012-04 -1.601631
2012-05 0.778033
Freq: M, dtype: float64
In [118]: ps.to_timestamp()
Out[118]:
2012-01-01 -1.475051
2012-02-01 0.722570
2012-03-01 -0.322646
2012-04-01 -1.601631
2012-05-01 0.778033
Freq: MS, dtype: float64
```

Converting between period and timestamp enables some convenient arithmetic functions to be used. In the following example, we convert a quarterly frequency with year ending in November to 9am of the end of the month following the quarter end:

```
In [119]: prng = pd.period_range("1990Q1", "2000Q4", freq="Q-NOV")
In [120]: ts = pd.Series(np.random.randn(len(prng)), prng)
In [121]: ts.index = (prng.asfreq("M", "e") + 1).asfreq("H", "s") + 9
In [122]: ts.head()
Out[122]:
1990-03-01 09:00    -0.289342
1990-06-01 09:00    0.233141
1990-09-01 09:00    -0.223540
1990-12-01 09:00    0.542054
1991-03-01 09:00    -0.688585
Freq: H, dtype: float64
```

Categoricals

pandas can include categorical data in a <u>DataFrame</u>. For full docs, see the <u>categorical introduction</u> and the <u>API documentation</u>.

Convert the raw grades to a categorical data type.

Rename the categories to more meaningful names (assigning to <u>Series.cat.categories()</u> is in place!).

```
In [126]: df["grade"].cat.categories = ["very good", "good", "very bad"]
```

Reorder the categories and simultaneously add the missing categories (methods under <u>Series.cat()</u> return a new <u>Series</u> by default).

```
In [127]: df["grade"] = df["grade"].cat.set_categories(
              ["very bad", "bad", "medium", "good", "very good"]
   ....:)
   • • • • • •
In [128]: df["grade"]
Out[128]:
    very good
1
          good
2
          good
3
    very good
    very good
     very bad
Name: grade, dtype: category
Categories (5, object): ['very bad', 'bad', 'medium', 'good', 'very good']
```

Sorting is per order in the categories, not lexical order.

```
In [129]: df.sort_values(by="grade")
Out[129]:
  id raw_grade
                grade
       e very bad
 6
1 2
         b
                 good
2 3
         b
                 good
0 1
          a very good
3 4
          a very good
4 5
          a very good
```

Grouping by a categorical column also shows empty categories.

```
In [130]: df.groupby("grade").size()
Out[130]:
grade
very bad     1
bad      0
medium      0
good      2
very good     3
dtype: int64
```

Plotting

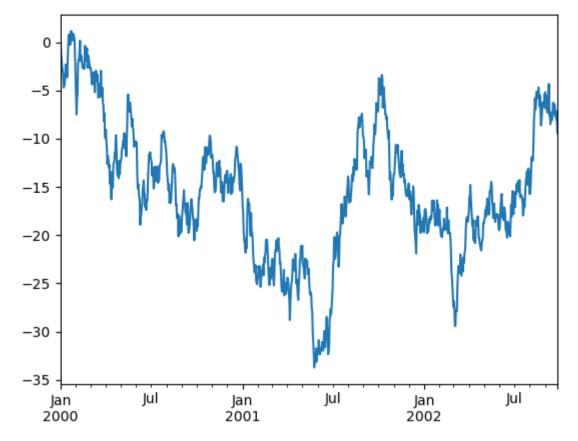
See the <u>Plotting</u> docs.

We use the standard convention for referencing the matplotlib API:

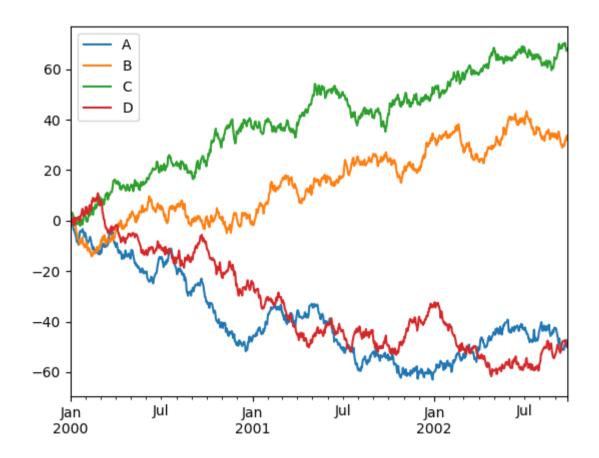
```
In [131]: import matplotlib.pyplot as plt
In [132]: plt.close("all")
```

The **close()** method is used to <u>close</u> a figure window.

```
In [133]: ts = pd.Series(np.random.randn(1000), index=pd.date_range("1/1/2000", periods=1000))
In [134]: ts = ts.cumsum()
In [135]: ts.plot();
```



On a DataFrame, the plot() method is a convenience to plot all of the columns with labels:



Getting data in/out

CSV

Writing to a csv file.

```
In [141]: df.to_csv("foo.csv")
```

Reading from a csv file.

```
In [142]: pd.read_csv("foo.csv")
Out[142]:
    Unnamed: 0
                     Α
                               В
                                         C
    2000-01-01 0.350262 0.843315 1.798556 0.782234
    2000-01-02 -0.586873 0.034907 1.923792 -0.562651
1
    2000-01-03 -1.245477 -0.963406 2.269575 -1.612566
    2000-01-04 -0.252830 -0.498066 3.176886 -1.275581
3
    2000-01-05 -1.044057 0.118042 2.768571 0.386039
995 2002-09-22 -48.017654 31.474551 69.146374 -47.541670
996 2002-09-23 -47.207912 32.627390 68.505254 -48.828331
997 2002-09-24 -48.907133 31.990402 67.310924 -49.391051
998 2002-09-25 -50.146062 33.716770 67.717434 -49.037577
999 2002-09-26 -49.724318 33.479952 68.108014 -48.822030
[1000 rows x \ 5 columns]
```

HDF5

Reading and writing to **HDFStores**.

Writing to a HDF5 Store.

```
In [143]: df.to_hdf("foo.h5", "df")
```

Reading from a HDF5 Store.

```
In [144]: pd.read_hdf("foo.h5", "df")
Out[144]:
                            В
                                      C
2000-01-01 0.350262 0.843315 1.798556 0.782234
2000-01-02 -0.586873 0.034907
                                1.923792 -0.562651
2000-01-03 -1.245477 -0.963406 2.269575 -1.612566
2000-01-04 -0.252830 -0.498066 3.176886 -1.275581
2000-01-05 -1.044057 0.118042 2.768571 0.386039
               . . .
                          . . .
                                     . . .
2002-09-22 -48.017654 31.474551 69.146374 -47.541670
2002-09-23 -47.207912 32.627390 68.505254 -48.828331
2002-09-24 -48.907133 31.990402 67.310924 -49.391051
2002-09-25 -50.146062 33.716770 67.717434 -49.037577
2002-09-26 -49.724318 33.479952 68.108014 -48.822030
[1000 rows x + 4 columns]
```

Excel

Reading and writing to MS Excel.

Writing to an excel file.

```
In [145]: df.to_excel("foo.xlsx", sheet_name="Sheet1")
```

Reading from an excel file.

```
In [146]: pd.read_excel("foo.xlsx", "Sheet1", index_col=None, na_values=["NA"])
   Unnamed: 0
                                          C
   2000-01-01 0.350262 0.843315 1.798556
                                             0.782234
1 2000-01-02 -0.586873 0.034907 1.923792 -0.562651
2 2000-01-03 -1.245477 -0.963406 2.269575 -1.612566
3 2000-01-04 -0.252830 -0.498066 3.176886 -1.275581
4 2000-01-05 -1.044057 0.118042 2.768571 0.386039
               . . .
                                    . . .
995 2002-09-22 -48.017654 31.474551 69.146374 -47.541670
996 2002-09-23 -47.207912 32.627390 68.505254 -48.828331
997 2002-09-24 -48.907133 31.990402 67.310924 -49.391051
998 2002-09-25 -50.146062 33.716770 67.717434 -49.037577
999 2002-09-26 -49.724318 33.479952 68.108014 -48.822030
[1000 rows x 5 columns]
```

Gotchas

If you are attempting to perform an operation you might see an exception like:

```
>>> if pd.Series([False, True, False]):
... print("I was true")
Traceback
...
ValueError: The truth value of an array is ambiguous. Use a.empty, a.any() or a.all().
```

See <u>Comparisons</u> for an explanation and what to do.

See Gotchas as well.

<< User Guide

Intro to data structures >>

© Copyright 2008-2021, the pandas development team. Created using <u>Sphinx</u> 4.1.2.