

Python Cheat Sheet

Pandas Basics



Pandas

The **Pandas** library is built on NumPy and provides easy-to-use **data structures** and **data analysis** tools for the Python programming language.



Use the following import convention:

```
>>> import pandas as pd
```

Pandas Data Structures

Series

A **one-dimensional** labeled array capable of holding any data type

a	3
b	-5
c	5
d	7

Index → 4

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

DataFrame

Columns

	Country	Capital	Population
0	Belgium	Brussels	11190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

A **two-dimensional** labeled data structure with columns of potentially different types

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [11190846, 1303171035, 207847528]}
>>> df = pd.DataFrame(data,
   columns=['Country', 'Capital', 'Population'])
```

I/O

Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)>>> from sqlalchemy import create_engine >>>
df.to_csv('myDataFrame.csv')>>> engine = create_engine('sqlite:///:memory:')
```

Read and Write to Excel

```
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_excel('file.xlsx')
```

```
>>> pd.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
Read multiple sheets from the same file read_sql() is a convenience wrapper around read_sql_table() and read_sql_query()
```

```
>>> xlsx = pd.ExcelFile('file.xls')
>>> df = pd.read_excel(xlsx, 'Sheet1')
```

Asking For

```
>>> help(pd.Series.loc)
```

Help Selection

Getting

```
>>> s['b']
-5
>>> df[1:]
   Country Capital Population
1  India     New Delhi  1303171035
2  Brazil    Brasilia  207847528
```

Also see NumPy Arrays

Get one element
Get subset of a DataFrame

Dropping

```
>>> s.drop(['a', 'c']) Drop values from rows (axis=0)
>>> df.drop('Country', axis=1) Drop values from columns (axis=1)
```

Sort & Rank

```
>>> df.sort_index() Sort by labels along an axis
>>> df.sort_values(by='Country') Sort by the values along an axis
>>> df.rank() Assign ranks to entries
```

Retrieving Series/DataFrame Information

Basic Information

```
>>> df.shape (rows, columns)
>>> df.index Describe index
>>> df.columns Describe DataFrame columns Info on DataFrame
>>> df.info() Number of non-NA values
```

Summary

```
>>> df.sum() Cumulative sum
>>> df.cumsum() Minimum/maximum
>>> df.idxmin() / df.idxmax() Summary statistics
>>> df.mean() Mean of values
>>> df.median() Median of values
```

Sum of values of values values Minimum/Maximum index value

Applying Functions

```
>>> f = lambda x: x**2
>>> df.apply(f) Apply function
>>> df.applymap(f) Apply function element-wise
```

Data Alignment

Internal Data Alignment

NA values are introduced in the indices that don't overlap:

```
>>> s3 = pd.Series([7, -2, 3], index=['a', 'c', 'd'])
>>> s + s3
   a    10.0
   b      NaN
   c    5.0
   d    7.0
```

Arithmetic Operations with Fill Methods

You can also do the internal data alignment yourself with the help of the fill methods:

```
>>> pd.to_sql('myDF', engine)
```

```
>>> s.add(s3, fill_value=0) a 10.0
   b -5.0   c
5.0
   d 7.0
>>>           s.sub(s3,
fill_value=2)>>>
s.div(s3,fill_value=4)
>>>           s.mul(s3,
fill_value=3)
```

