

# Python For Data Science Cheat Sheet

# Pandas Basics

## Learn Python for Data Science Interactively

### 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

Index →	a	3
	b	-5
	c	7
	d	4

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

#### DataFrame

Columns →	Country	Capital	Population
Index →	0	Belgium	Brussels
	1	India	New Delhi
	2	Brazil	Brasilia

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'])
```

### Asking For Help

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

### Selection

Also see NumPy Arrays

#### Getting

<pre>&gt;&gt;&gt; s['b'] -5 &gt;&gt;&gt; df[1:]    Country  Capital  Population 1  India    New Delhi  1303171035 2  Brazil   Brasilia   207847528</pre>	<p>Get one element</p> <p>Get subset of a DataFrame</p>
--	---

#### Selecting, Boolean Indexing & Setting

<p><b>By Position</b></p> <pre>&gt;&gt;&gt; df.iloc[[0],[0]] 'Belgium' column &gt;&gt;&gt; df.iat([0],[0]) 'Belgium'<p><b>By Label</b></p><pre>&gt;&gt;&gt; df.loc[[0], ['Country']] 'Belgium' column labels &gt;&gt;&gt; df.at([0], ['Country']) 'Belgium'<p><b>By Label/Position</b></p><pre>&gt;&gt;&gt; df.ix[2] Country      Brazil Capital    Brasilia Population  207847528 &gt;&gt;&gt; df.ix[:, 'Capital'] 0 Brussels 1 New Delhi 2 Brasilia &gt;&gt;&gt; df.ix[1, 'Capital'] 'New Delhi'<p><b>Boolean Indexing</b></p><pre>&gt;&gt;&gt; s[~(s &gt; 1)] &gt;&gt;&gt; s[(s &lt; -1)   (s &gt; 2)] &gt;&gt;&gt; df[df['Population']&gt;1200000000]<p><b>Setting</b></p><pre>&gt;&gt;&gt; s['a'] = 6</pre></pre></pre></pre></pre>	<p>Select single value by row &amp;</p> <p>Select single value by row &amp;</p> <p>Select single row of subset of rows</p> <p>Select a single column of subset of columns</p> <p>Select rows and columns</p> <p>Series s where value is not &gt;1 s where value is &lt;-1 or &gt;2 Use filter to adjust DataFrame</p> <p>Set index a of Series s to 6</p>
--	---

### Dropping

<pre>&gt;&gt;&gt; s.drop(['a', 'c']) &gt;&gt;&gt; df.drop('Country', axis=1)</pre>	<p>Drop values from rows (axis=0)</p> <p>Drop values from columns (axis=1)</p>
--	--

### Sort & Rank

<pre>&gt;&gt;&gt; df.sort_index() &gt;&gt;&gt; df.sort_values(by='Country') &gt;&gt;&gt; df.rank()</pre>	<p>Sort by labels along an axis</p> <p>Sort by the values along an axis</p> <p>Assign ranks to entries</p>
--	--

### Retrieving Series/DataFrame Information

#### Basic Information

<pre>&gt;&gt;&gt; df.shape &gt;&gt;&gt; df.index &gt;&gt;&gt; df.columns &gt;&gt;&gt; df.info() &gt;&gt;&gt; df.count()</pre>	<p>(rows, columns)</p> <p>Describe index</p> <p>Describe DataFrame columns</p> <p>Info on DataFrame</p> <p>Number of non-NA values</p>
---	--

#### Summary

<pre>&gt;&gt;&gt; df.sum() &gt;&gt;&gt; df.cumsum() &gt;&gt;&gt; df.min()/df.max() &gt;&gt;&gt; df.idxmin()/df.idxmax() &gt;&gt;&gt; df.describe() &gt;&gt;&gt; df.mean() &gt;&gt;&gt; df.median()</pre>	<p>Sum of values</p> <p>Cummulative sum of values</p> <p>Minimum/maximum values</p> <p>Minimum/Maximum index value</p> <p>Summary statistics</p> <p>Mean of values</p> <p>Median of values</p>
--	--

### Applying Functions

<pre>&gt;&gt;&gt; f = lambda x: x*2 &gt;&gt;&gt; df.apply(f) &gt;&gt;&gt; df.applymap(f)</pre>	<p>Apply function</p> <p>Apply function element-wise</p>
--	--

### I/O

#### Read and Write to CSV

```
>>> pd.read_csv('file.csv', header=None, nrows=5)
>>> df.to_csv('myDataFrame.csv')
```

#### Read and Write to Excel

```
>>> pd.read_excel('file.xlsx')
>>> df.to_excel('dir/myDataFrame.xlsx', sheet_name='Sheet1')
```

Read multiple sheets from the same file

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

#### Read and Write to SQL Query or Database Table

```
>>> from sqlalchemy import create_engine
>>> engine = create_engine('sqlite:///memory:')
>>> pd.read_sql("SELECT * FROM my_table;", engine)
>>> pd.read_sql_table('my_table', engine)
>>> pd.read_sql_query("SELECT * FROM my_table;", engine)
```

read\_sql() is a convenience wrapper around read\_sql\_table() and read\_sql\_query()

```
>>> df.to_sql('myDf', engine)
```

### 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
```

#### Summary

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

```
>>> 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)
```