

Python For Data Science 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	-
c	5
d	7
	4

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

DataFrame

Columns

	Country	Capital	Population
0	Belgium	Brussels	1190846
1	India	New Delhi	1303171035
2	Brazil	Brasilia	207847528

```
>>> data = {'Country': ['Belgium', 'India', 'Brazil'],
   'Capital': ['Brussels', 'New Delhi', 'Brasilia'],
   'Population': [1190846, 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_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')
>>> pd.to_sql('myDF', engine)
```

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

Selecting, Boolean Indexing & Setting

By Position

```
>>> df.iat[0, 0]
'Belgium'
>>> df.iat[0, 0]
'Belgium'
```

By Label

```
>>> df.loc[0, 'Country']
'Belgium'
>>> df.at[0, 'Country']
'Belgium'
```

By Label/Position

```
>>> df.ix[2]
Country Brazil
Capital Brasilia
Population 207847528
>>> df.ix[:, 'Capital']
0 Brussels
1 New Delhi
2 Brasilia
>>> df.ix[1, 'Capital']
'New Delhi'
```

Boolean Indexing

```
>>> s[s > 1]
>>> s[(s < -1) | (s > 2)]
>>> df[df['Population'] > 1200000000]
```

Setting

```
>>> s['a'] = 6
```

Get one element

Get subset of a DataFrame

Select single row of subset of rows

Series s where value is not >1
where value is <-1 or >2

Use filter to adjust DataFrame

Set index

a of Series s to 6

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()
>>> df.sort_values(by='Country')
>>> df.rank()
```

Sort by labels along an axis
Sort by the values along an axis
Assign ranks to entries

Retrieving Series/DataFrame Information

Basic Information

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

Summary

>>> df.sum()	Sum of values
>>> df.cumsum()	Cumulative sum of values
>>> df.min()	Minimum/maximum values
>>> df.idxmin()	Minimum/Maximum index value
>>> df.describe()	Summary statistics
>>> df.mean()	Mean of values
>>> df.median()	Median of values

Applying Functions

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

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
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:

```
>>> s.add(s3, fill_value=0)
a 10.0
b -5.0
c 5.0
d 7.0
>>> s.sub(s3, fill_value=2)
a 8.0
b -7.0
c 3.0
d 5.0
>>> s.div(s3, fill_value=4)
a 2.5
b -1.0
c 1.25
d 1.75
>>> s.mul(s3, fill_value=3)
a 30.0
b -15.0
c 15.0
d 21.0
```

