

“Advanced Python Programming for Everybody”

Instructor: Ernest Bonat, Ph.D.

Senior Software Engineer

Senior Data Scientist

ebonat@15itresources.com

Cell: 503.730.4556

Module 5 Source Code

https://github.com/ebonat/intel_module_5

Module 5. “Python Data Ecosystem for Data Science Projects – Part 1”

What do you really need to know to become a Data Scientist?

- **Probability and Statistics** (undergraduate level)
- **Python Programming Language** (good level!)
- **Python Data Ecosystem** (good level!):
 1. **NumPy** – fundamental package for scientific computing (Numerical Python - <http://www.numpy.org/>)
 2. **pandas** – provides easy-to-use and high-performance data structures (<https://pandas.pydata.org/>)
 3. **SciPy** - Python-based ecosystem of open-source software for mathematics, science, and engineering (<https://www.scipy.org/>)

4. **scikit-learn Machine Learning** – a simple and efficient tool for data mining and data analysis (<http://scikit-learn.org/>)
5. **matplotlib** – a 2D plotting library which produces publication quality figures in a variety of hard copy formats and interactive environments across platforms (<https://matplotlib.org/>)
6. **seaborn** - statistical data visualization (<https://seaborn.pydata.org/>)
7. **scikit-image** – a collection of algorithms for image processing (<http://scikit-image.org/>)

Best way to learn any of them? – **have a data project to do!**

Data Science Two Main Tasks:

1	Data Cleansing (in Python is done with pandas) or Data Preprocessing. Require a “Data Cleansing Requirements Document”	60% - 70% work
2	Data Analytics	40% - 30% work

Data cleansing very important task, so be careful with “**Garbage IN – Garbage OUT**”

Why pandas?

- Heterogeneous data types
- Easy, fast missing data handling
- Easier to write generic code
- Labeled data (numpy mostly assumes index == label)
- Relational data

Install pandas

```
conda install pandas or conda update pandas
```

pandas - an open source library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.

(<http://pandas.pydata.org>)

Cheat Sheet

https://github.com/pandas-dev/pandas/blob/master/doc/cheatsheet/Pandas_Cheat_Sheet.pdf

PDF Documentation File

<http://pandas.pydata.org/pandas-docs/stable/pandas.pdf>

Created by: **Wes McKinney**, now maintained by Jeff Reback and many others

O'Reilly Media Book: **Python for Data Analysis, 2nd Edition**

Data Wrangling with Pandas, NumPy, and IPython

(<http://shop.oreilly.com/product/0636920050896.do>)

Pandas Cheat Sheet - Python for Data Science

<https://www.dataquest.io/blog/pandas-cheat-sheet/>

Two main imports:

```
import numpy as np
```

```
import pandas as pd
```

Definition

df	Any pandas DataFrame object
s	Any pandas Series object

Importing Data

pd.read_csv(filename) (very slow for big files – use Pandas on Ray https://rise.cs.berkeley.edu/blog/pandas-on-ray/)	From a CSV file
pd.read_table(filename)	From a delimited text file (like TSV)
pd.read_excel(filename)	From an Excel file
pd.read_sql(query, connection_object)	Read from a SQL table/database
pd.read_json(json_string)	Read from a JSON formatted string, URL or file.
pd.read_html(url)	Parses an html URL, string or file

	and extracts tables to a list of dataframes
<code>pd.read_clipboard()</code>	Takes the contents of your clipboard and passes it to <code>read_table()</code>
<code>pd.DataFrame(dict)</code>	From a dict, keys for columns names, values for data as lists

Exporting Data

<code>df.to_csv(filename)</code>	Write to a CSV file
<code>df.to_excel(filename)</code>	Write to an Excel file
<code>df.to_sql(table_name, connection_object)</code>	Write to a SQL table
<code>df.to_json(filename)</code>	Write to a file in JSON format

Create Test Objects

<code>pd.DataFrame(np.random.rand(20,5))</code>	5 columns and 20 rows of random floats
<code>pd.Series(my_list)</code>	Create a series from an iterable <code>my_list</code>
<code>df.index = pd.date_range('1900/1/30', periods=df.shape[0])</code>	Add a date index

Viewing/Inspecting Data

<code>df.head(n)</code>	First n rows of the DataFrame
<code>df.tail(n)</code>	Last n rows of the DataFrame
<code>df.shape()</code>	Number of rows and columns
<code>df.info()</code>	Index, Datatype and Memory information
<code>df.describe()</code>	Summary statistics for numerical columns
<code>s.value_counts(dropna=False)</code>	View unique values and counts

<code>df.apply(pd.Series.value_counts)</code>	Unique values and counts for all columns
---	--

Selection

<code>df[col]</code>	Return column with label col as Series
<code>df[[col1, col2]]</code>	Return Columns as a new DataFrame
<code>s.iloc[0]</code>	Selection by position
<code>s.loc['index_one']</code>	Selection by index
<code>df.iloc[0,:]</code>	First row
<code>df.iloc[0,0]</code>	First element of first column

Data Cleaning

<code>df.columns = ['a','b','c']</code>	Rename columns
<code>pd.isnull()</code>	Checks for null Values, Returns Boolean Array
<code>pd.notnull()</code>	Opposite of <code>pd.isnull()</code>
<code>df.dropna()</code>	Drop all rows that contain null values
<code>df.dropna(axis=1)</code>	Drop all columns that contain null values
<code>df.dropna(axis=1,thresh=n)</code>	Drop all rows have have less than n non null values
<code>df.fillna(x)</code>	Replace all null values with x
<code>s.fillna(s.mean())</code>	Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)
<code>s.astype(float)</code>	Convert the datatype of the series to float
<code>s.replace(1,'one')</code>	Replace all values equal to 1 with 'one'
<code>s.replace([1,3],['one','three'])</code>	Replace all 1 with 'one' and 3 with 'three'
<code>df.rename(columns=lambda x: x + 1)</code>	Mass renaming of columns
<code>df.rename(columns={'old_name':</code>	Selective renaming

<code>'new_name'})</code>	
<code>df.set_index('column_one')</code>	Change the index
<code>df.rename(index=lambda x: x + 1)</code>	Mass renaming of index

Joins

<code>df1.append(df2)</code>	Add the rows in df1 to the end of df2 (columns should be identical)
<code>df.concat([df1, df2],axis=1)</code>	Add the columns in df1 to the end of df2 (rows should be identical)
<code>df1.join(df2,on=col1,how='inner')</code>	SQL-style join the columns in df1 with the columns on df2 where the rows for col have identical values. how can be one of 'left', 'right', 'outer', 'inner'

Descriptive Statistics

(These can all be applied to a series as well)

<code>df.describe()</code>	Summary statistics for numerical columns
<code>df.mean()</code>	Return the mean of all columns
<code>df.corr()</code>	Finds the correlation between columns in a DataFrame
<code>df.count()</code>	Counts the number of non-null values in each DataFrame column
<code>df.max()</code>	Finds the highest value in each column
<code>df.min()</code>	Finds the lowest value in each column
<code>df.median()</code>	Finds the median of each column
<code>df.std()</code>	Finds the standard deviation of each column

pandas Data Structures Objects

1. Series

2.DataFrame

3.Panel (not very used today!)

Series

A one-dimensional labeled array capable of holding any data type (integers, strings, floating point numbers, Python objects, etc.). The axis labels are collectively referred to as the **index**.

```
s = pd.Series(data, index=index)
```

Where: data can be: Python dictionary, ndarray (n-dimensional array or any scalar value (like 10))

Example:

```
s = pd.Series(np.random.randn(5))  
print(s)
```

Result:

```
0    0.3674  
1   -0.8230  
2   -1.0295  
3   -1.0523  
4   -0.8502  
dtype: float64
```

DataFrame

A 2-dimensional labeled data structure with rows and columns of potentially different types (similar to Microsoft Excel spreadsheet or SQL database table)

```
df = pd.DataFrame(data, ...)
```

DataFrame accepts many different kinds of input:

- Dictionary of 1D ndarrays, lists, dicts, or Series
- 2-D numpy.ndarray
- Structured or record ndarray
- A Series
- Another DataFrame

Example:

```
dictionary = {"one" : [1., 2., 3., 4.], "two" : [4., 3., 2., 1.]}  
df = pd.DataFrame(dictionary)
```

Result:

```
   one two  
0  1.0  4.0  
1  2.0  3.0  
2  3.0  2.0  
3  4.0  1.0
```

Indexing / Selection

The basics of indexing are as follows:

Operation	Syntax	Result
Select column	<code>df[col]</code>	Series
Select row by label	<code>df.loc[label]</code>	Series
Select row by integer location	<code>df.iloc[loc]</code>	Series
Slice rows	<code>df[5:10]</code>	DataFrame
Select rows by boolean vector	<code>df[bool_vec]</code>	DataFrame

Panel

A 3-dimensional labeled data structure. It's less-used today!

Missing Data

Missing Data is define as Non-available (NA), null or “not present for whatever reason”

pandas uses “NaN” (Non-a-Number) or “nan” internally for simplicity and performance reasons

In CSV file:

one	two	three	four	five	timestamp
	2.1	3.1	bar	1	
	2.3	3.2		0	1/1/2017
1.3		3.3	bar	1	2/1/2017
1.4	2.4		bar		3/1/2017
1.5	2.5	3.5	bar	0	

In pandas DataFrame we'll have:

```
=====
```

```
..  one  two  three  four  five  timestamp
```

```
=====
```

```
0 nan    2.1    3.1 bar    1    nan
```

```
1 nan    2.3    3.2 nan    0    1/1/2017
```

2	1.3	nan	3.3	bar	1	2/1/2017
---	-----	-----	-----	-----	---	----------

3	1.4	2.4	nan	bar	nan	3/1/2017
---	-----	-----	-----	-----	-----	----------

4	1.5	2.5	3.5	bar	0	nan
---	-----	-----	-----	-----	---	-----

==== =====

pandas code examples!