

“Advanced Python Programming for Machine Learning Projects”

Instructor: Ernest Bonat, Ph.D.

Senior Software Engineer

Senior Data Scientist

ebonat@15itresources.com

Cell: 503.730.4556

2. Data Visualization and Pre-processing

GitHub: https://github.com/ebonat/intel_session_2

Data Science Two Main Tasks:

1	Data Pre-processing (Cleansing)	60% - 70% work
2	Data Analytics	40% - 30% work

Data Pre-processing very important task. Be careful with “Garbage IN – Garbage OUT”

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>

Two main imports:

```
import numpy as np  
import pandas as pd
```

Definition

df	Any pandas DataFrame object
ds	Any pandas (Data) Series object

Importing Data

pd.read_csv(filename) – slow with thousands of rows!	From a CSV file
pd.read_table(filename)	From a delimited text file (like

	TSV)
<code>pd.read_excel(filename)</code>	From an Excel file
<code>pd.read_sql(query, connection_object)</code>	Read from a SQL table/database
<code>pd.read_json(json_string)</code>	Read from a JSON formatted string, URL or file.
<code>pd.read_html(url)</code>	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 column names, values for data as lists

<code>pd.read_hdf(path_or_buf= filename, key="filename_key")</code> - very fast!	From a HDF5 file
---	------------------

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
<code>pd.to_hdf(path_or_buf= filename, key="filename_key")</code>	Write the contained data to an HDF5 file using HDFStore.

Viewing/Inspecting Data

df.head(n)	First n rows of the DataFrame
df.tail(n)	Last n rows of the DataFrame
df.shape()	Number of rows and columns
df.info()	Index, Datatype and Memory information
df.describe()	Summary statistics for numerical columns
ds.value_counts(dropna=False)	View unique values and counts
df.apply(pd.Series.value_counts)	Unique values and counts for all columns

Selection

df[col]	Return column with label col as Series
df[[col1, col2]]	Return Columns as a new DataFrame
ds.iloc[0]	Selection by position
ds.loc['index_one']	Selection by index
df.iloc[0,:]	First row
df.iloc[0,0]	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>ds.fillna(s.mean())</code>	Replace all null values with the mean (mean can be replaced with almost any function from the statistics section)

<code>ds.astype(float)</code>	Convert the datatype of the series to float
<code>ds.replace(1, 'one')</code>	Replace all values equal to 1 with 'one'
<code>ds.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': 'new_name'})</code>	Selective renaming
<code>df.set_index('column_one')</code>	Change the index
<code>df.rename(index=lambda x: x + 1)</code>	Mass renaming of index

Descriptive Statistics

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

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

Why pandas?

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

pandas Data Structures Objects

1.Series

2.DataFrame

3.Panel

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**.

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

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

Example:

```
ds = 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:

```
====  =====  =====  =====  =====  =====
```

```
..  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

==== =====

Beginning Steps:

- 1. Organize Input and Output Data Files Path Name**
- 2. Import Data File to pandas DataFrame**
- 3. Get Number of Rows and Columns**
- 4. Get Index, Datatype and Memory Information**
- 5. Remove Duplicates Rows**
- 6. Fill Nan Values (Mean, Median, Defaults, etc.)**
- 7. Remove Rows by Row/Column Conditions**
- 8. Replace Values by Row/Column Conditions**