**Pandas and Numpy Quick Reference Sheet**
**Version 1.0.3 by Joff Thyer**
**Rivergum Security LLC**

| Concept | Description | Additional Information |
|---|---|---|
| import pandas as pd | conventional way to import into Python | Pandas and numpy tend to go together in pairs |
| import numpy as np | conventional way to import into Python | |
| Pandas "Series" object | a one dimensional labeled array | |
| Pandas "DataFrame" object | a two dimensional data structure that holds a table | |
| s = pd.Series([1,2,np.nan,4]) | creating a Pandas series | |
| dates = pd.date_range("20130101", periods=6) | create a Pandas date range with 6 periods (defaults to days) | |
| df = pd.DataFrame(np.random.randn(2, 4), index=dates, columns=list("ABCD")) | generate a dataframe with index of date range, and 4 columns. | |
| np.random.randn(6, 4) | creates a 6 x 4 random floating point array from standard normal distribution | |
| `df2 = pd.DataFrame({ "A": 1.0, "B": pd.Timestamp("20130102"), "C": pd.Series(1, index=list(range(4)), dtype="float32"), "D": np.array([3] * 4, dtype="int32"), "E": pd.Categorical(["test", "train", "test", "train"]), "F": "foo", })` | generates a dataframe from a dictionary whereby the column names are the keys in the dictionary | |
| DataFrame.head() | view the first few rows of a dataframe | |
| DataFrame.tail() | view the last few rows of a dataframe | |
| DataFrame.to_numpy() | return the numpy array of a dataframe without index of column names | |
| df.describe() DataFrame.describe() | yield basic statistical information about a dataframe. | |
| df.T | transpose the data (matrix) | |
| df.sort_index(axis=1, ascending=False) | sort by axis 1 (rows). Axis=0 would be sort by columns. | |
| df.sort_values(by="B") | sort by values in column "B" | |
| df[0:2] | select two rows by slicing | |
| df["2013-01-01":"2013-01-02"] | select rows by slicing on date index | |
| df["A"] | select data by column name | |
| df.loc[dates[0]] | select data by label | |
| df.loc['2013-01-02':'2013-01-04', ['A', 'D']] | slice by dates and select specific columns | |
| df.at[dates[0], 'C'] | fast access to a single scalar in column C | |
| DataFrame.iloc() | select by position | |
| DataFrame.iat() | select by position | |
| df.iloc[1:3, :] | slicing rows explicity | |
| df.iloc[:, 1:3] | slicing columns explicity | |
| df.mean() | calculate the mean of column data | |
| df.mean(axis=1) | calculate the mean of rows | |
| DataFrame.agg() | aggregate user a user defined function over a specified axis | |
| DataFrame.transform() | call a user defined function across the whole dataframe | |
| df.value_counts() | perform a frequency count on the dataframe | |