

Stata to Python Equivalents

Contents

- [1 Intro/Note on Notation](#)
- [2 Input/Output](#)
- [3 Sample Selection](#)
- [4 Data Info and Summary Statistics](#)
- [5 Variable Manipulation](#)
- [6 Panel Data](#)
- [7 Merging and Joining](#)
- [8 Reshape](#)
- [9 Econometrics](#)
- [10 Plotting](#)
- [11 Other differences](#)

Special thanks to John Coglianese for feedback and for supplying the list of "vital" Stata commands. Feedback and requests for additions to the list are always welcome!

The official Pandas documentation includes a "[Comparison with Stata](https://pandas.pydata.org/pandas-docs/stable/comparison_with_stata.html)" (https://pandas.pydata.org/pandas-docs/stable/comparison_with_stata.html) page which is another great resource.

[1 Intro/Note on Notation](#)

Coding in Python is a little different than coding in Stata.

In Stata, you have one dataset in memory. The dataset is a matrix where each column is a "variable" with a unique name and each row has a number (the special variable `_n`). Everything in Stata is built around this paradigm.

Python is a general purpose programming language where a "variable" is not a column of data. Variables can be anything, a single number, a matrix, a list, a string, etc. The Pandas package implements a kind of variable called a DataFrame that acts a lot like the single dataset in Stata. It is a matrix where each column and each row has a name. The key distinction in Python is that a DataFrame is itself a variable and you can work with any number of DataFrames at one time. You can think of each column in a DataFrame as a variable just like in Stata, except that when you reference a column, you also have to specify the DataFrame.

The Stata-to-Python translations below are written assuming that you have a single DataFrame called `df`. Placeholders like `<varname>` and `<datafile>` show where user-specified values go in each language. Note that in many cases, `<varname>` will be simple text in Stata (e.g., `avg_income`) while in Python it will be a string (`'avg_income'`). If you were to write `df[avg_income]` without quotes, Python would go looking for a variable--a list, a number, a string--that's been defined somewhere else. Because of this, `<varlist>` in Python represents a list of variable names: `['educ', 'income', 2017]`.

2 Input/Output

Stata	Python
log using <file>	Python doesn't display results automatically like Stata. You have to explicitly call the <code>print</code> function. Using a Jupyter notebook is the closest equivalent.
help <command>	<code>help(<command>)</code> OR <code><command>?</code> in IPython (as in <code>pd.read_stata?</code>)
cd <directory>	<code>import os</code> <code>os.chdir('<directory>')</code> but this is bad practice. Better practice is to use full pathnames whenever possible.
use <dtafile>	<code>import pandas as pd</code> <code>df = pd.read_stata('<dtafile>')</code>
use <varlist> using <dtafile>	<code>df = pd.read_stata('<dtafile>', columns=<varlist>)</code>
import excel using <excelfile>	<code>df = pd.read_excel('<excelfile>')</code>
import delimited using <csvfile>	<code>df = pd.read_csv('<csvfile>')</code>
save <filename>, replace	<code>df.to_stata('<filename>')</code> OR <code>df.to_pickle('<filename>')</code> for Python-native file type.
outsheet using <csv_name>, comma	<code>df.to_csv('<csv_name>')</code>
export excel using <excel_name>	<code>df.to_excel('<excel_name>')</code>

3 Sample Selection

Stata	Python
keep if <condition>	<code>df = df[<condition>]</code>
drop if <condition>	<code>df = df[~(<condition>)]</code>
keep <var>	<code>df = df[<var>]</code>
keep varstem*	<code>df = df.filter(like='varstem*')</code>
drop <var>	<code>del df[<var>]</code> OR <code>df = df.drop(<var>, axis=1)</code>
drop varstem*	<code>df =</code> <code>df.drop(df.filter(like='varstem*').columns, axis=1)</code>

4 Data Info and Summary Statistics

Stata	Python
<code>describe</code>	<code>df.info()</code> OR <code>df.dtypes</code> just to get data types. Note that Python does not have value labels like Stata does.
<code>describe <var></code>	<code>df[<var>].dtype</code>
<code>count</code>	<code>df.shape[0]</code> OR <code>len(df)</code> . Here <code>df.shape</code> returns a tuple with the length and width of the DataFrame.
<code>count if <condition></code>	<code>df[<condition>].shape[0]</code> OR <code>(<condition>).sum()</code> if the condition involves a DataFrame, e.g., <code>(df['age'] > 2).sum()</code>
<code>summ <var></code>	<code>df['<var>'].describe()</code>
<code>summ <var> if <condition></code>	<code>df[<condition>][<var>].describe()</code> OR <code>df.loc[<condition>, <var>].describe()</code>
<code>summ <var> [aw = <weight>]</code>	Right now you have to calculate weighted summary stats manually. There are also some tools available in the Statsmodels package.
<code>summ <var>, d</code>	<code>df[<var>].describe()</code> plus <code>df[<var>].quantile([.1, .25, .5, .75, .9])</code> or whatever other statistics you want.

5 Variable Manipulation

Stata	Python
<code>gen <newvar> = <expression></code>	<code>df[<newvar>] = <expression></code>
<code>gen <newvar> = <expression> if <condition></code>	<code>df.loc[<condition>, <newvar>] = <expression></code> . As with Stata, the rows of <code>df</code> that don't meet the condition will be missing (<code>numpy.nan</code>).
<code>replace <var> = <expression> if <condition></code>	<code>df.loc[<condition>, <var>] = <expression></code>
<code>rename <var> <newvar></code>	<code>df = df.rename(columns={<var>: <newvar>})</code> . You can also directly manipulate <code>df.columns</code> like a list: <code>df.columns = ['a', 'b', 'c']</code> .
<code>inlist(<var>, <val1>, <val2>)</code>	<code>df[<var>].isin((<val1>, <val2>))</code>

Stata	Python
<code>inrange(<var>, <val1>, <val2>)</code>	<code>df[<var>].between((<val1>, <val2>))</code>
<code>subinstr(<str>, " ", "_", .)</code>	<code>df[<var>].str.replace(' ', '_')</code>
<code>egen <newvar> = count(<var>)</code>	<code><newvar> = df[<var>].notnull().sum()</code> NOTE: For these egen commands, <newvar> is a full (constant) column in Stata, while it is a scalar in Python.
<code>egen <newvar> = group(<varlist>)</code>	<code><newvar> = econtools.group_id(df, cols=<varlist>)</code>
<code>egen <newvar> = max(<var>)</code>	<code><newvar> = df[<var>].max()</code>
<code>egen <newvar> = mean(<var>)</code>	<code><newvar> = df[<var>].mean()</code>
<code>egen <newvar> = total(<var>)</code>	<code><newvar> = df[<var>].sum()</code>
<code>egen <newvar> = <stat>(<var>), by(<groupvars>)</code>	<code>df[<newvar>] = df.groupby(<groupvars>)[<var>].transform('<stat>')</code>
<code>collapse (sd) <var> (median) <var> /// (max) <var> (min) <var>, /// by(<groupvars>)</code>	<code>df.groupby(<groupvars>)[<var>].agg(['std', 'median', 'min', 'max', 'sum'])</code>
<code>collapse (<stat>) <var> [iw = <weight>]</code>	Manually or maybe through Statsmodels tool.
<code>collapse (<stat>) <stat_vars>, by(<groupvars>)</code>	<code>df.groupby(<groupvars>)[<stat_vars>].<stat>()</code>
<code>recode <var> (1/5 = 1)</code>	N/A, see note below.
<code>recode <var> (1/5 = 1), gen(<newvar>)</code>	N/A.
<code>label var <var> <label></code>	N/A.
<code>label define <labelname> 1 <valuelabel></code>	N/A.
<code>label values <var> <labelname></code>	N/A.
<code>label list <labelname></code>	N/A.

Python doesn't have "labels" built into DataFrames like Stata does. However, you can use a dictionary to map data values to labels when necessary.

```
variable_labels = {
    1: "First Category",
    2: "Second Category",
    3: "Last Category",
}
```

6 Panel Data

There is no general equivalent to `tsset` in Python. However, you can accomplish most if not all of the same tasks using a DataFrame's index (the row's equivalent of columns.) In Stata, the "DataFrame" in memory always has the observation row number, denoted by the Stata built-in variable `_n`. In Python and Pandas, a DataFrame index can be anything (though you can also refer to rows by the row number; see `.loc` vs `iloc`). It can also be hierarchical with multiple levels. It is a much more general tool than `tsset`.

Stata	Python
<code>tsset <panelvar> <timevar></code>	<code>df = df.set_index([<panelvar>, <timevar>])</code>
<code>L.<var></code>	<code>df.shift()</code> NOTE: The index must be correctly sorted for <code>shift</code> to work the way you want it to.
<code>L2.<var></code>	<code>df.shift(2)</code>
<code>F.<var></code>	<code>df.shift(-1)</code>

6.1 Examples

```
In [1]: import numpy as np
```

```
In [2]: import pandas as pd
```

```
In [3]: df0 = pd.DataFrame({'var1': np.arange(6),
....:                      'id': [1, 1, 2, 2, 3, 3],
....:                      'period': [0, 1] * 3})
```

```
In [4]: print(df0)
   var1  id  period
0     0   1     0
1     1   1     1
2     2   2     0
3     3   2     1
4     4   3     0
5     5   3     1
```

```
In [5]: df = df0.set_index(['id', 'period'])
```

```
In [6]: print(df)
      var1
id period
1  0     0
   1     1
2  0     2
   1     3
3  0     4
   1     5
```

```
In [7]: df['var1_lag'] = df.groupby(level='id')['var1'].shift()
```

```
In [8]: print(df)
      var1  var1_lag
id period
1  0     0      NaN
   1     1     0.0
2  0     2      NaN
   1     3     2.0
3  0     4      NaN
   1     5     4.0
```

```
In [9]: df['var1_for'] = df.groupby(level='id')['var1'].shift(-1)
```

```
In [10]: print(df)
      var1  var1_lag  var1_for
id period
1  0     0      NaN     1.0
   1     1     0.0     NaN
2  0     2      NaN     3.0
   1     3     2.0     NaN
```

3	0	4	NaN	5.0
	1	5	4.0	NaN

7 Merging and Joining

Stata

```
append using <filename>

merge 1:1 <vars> using <filename>
```

Python

```
df_joint = df1.append(df2)

df_joint = df1.join(df2) if <vars> are the
DataFrames' indexes, or
df_joint = pd.merge(df1, df2, on=
<vars>) otherwise. Beware
that pd.merge will not keep the index of either
DataFrame.
NOTE: Merging in Python is like R, SQL, etc. Needs
more robust
explanation.
```

Merging with Pandas DataFrames does not require you to specify "many-to-one" or "one-to-many". Pandas will figure that out based on whether the variables you're merging on are unique or not. However, you can specify what sub-sample of the merge to keep using the keyword argument `how`, e.g., `df_joint = df1.join(df2, how='left')` is the default for `join` while `how='inner'` is the default for `pd.merge`.

Pandas how	Stata , keep()	Intuition
<code>how='left'</code>	<code>keep(1, 3)</code>	Keeps all observations in the "left" DataFrame.
<code>how='right'</code>	<code>keep(2, 3)</code>	Keeps all observations in the "right" DataFrame.
<code>how='inner'</code>	<code>keep(3)</code>	Keeps observations that are in both DataFrames.
<code>how='outer'</code>	<code>keep(1 2 3)</code>	Keeps all observations.

8 Reshape

Like with merging, reshaping a DataFrame in Python is a bit different because of the paradigm shift from the "only data table in memory" model of Stata to "a data table is just another object/variable" of Python. But this difference also makes reshaping a little easier in Python.

The most fundamental reshape commands in Python/Pandas are `stack` and `unstack`:

```
In [1]: import pandas as pd
```

```
In [2]: import numpy as np
```

```
In [3]: long = pd.DataFrame(np.arange(8),
....:                        columns=['some_variable'],
....:                        index=pd.MultiIndex.from_tuples(
....:                            [('a', 1), ('a', 2),
....:                             ('b', 1), ('b', 2),
....:                             ('c', 1), ('c', 2),
....:                             ('d', 1), ('d', 2)]))
```

```
In [4]: long.index.names=['unit_id', 'time']
```

```
In [5]: long.columns.name = 'varname'
```

```
In [6]: long
```

```
Out[6]:
```

varname		some_variable
unit_id	time	
a	1	0
	2	1
b	1	2
	2	3
c	1	4
	2	5
d	1	6
	2	7

```
In [7]: wide = long.unstack('time')
```

```
In [8]: wide
```

```
Out[8]:
```

varname		some_variable	
time		1	2
unit_id			
a		0	1
b		2	3
c		4	5
d		6	7

```
In [9]: long2 = wide.stack('time')
```

```
In [10]: long2
```

```
Out[10]:
```

varname		some_variable
unit_id	time	
a	1	0
	2	1
b	1	2
	2	3

c	1	4
	2	5
d	1	6
	2	7

Here Input 3 creates a DataFrame, Input 4 gives each of the index columns a name, and Input 5 names the columns. Coming from Stata, it's a little weird to think of the column names themselves having a "name", but the columns names are just an index like the row names are. It starts to make more sense when you realize columns don't have to be strings. They can be integers, like years or FIPS codes. In those cases, it makes a lot of sense to give the columns a name so you know what you're dealing with.

Input 6 does the reshaping using `unstack('time')`, which takes the index `'time'` and creates a new column for every unique value it has. Notice that the columns now have multiple levels, just like the index previously did. This is another good reason to label your index and columns. If you want to access either of those columns, you can do so as usual, using a tuple to differentiate between the two levels:

```
In [11]: wide[('some_variable', 1)]
Out[11]:
unit_id
a      0
b      2
c      4
d      6
Name: (some_variable, 1), dtype: int32
```

If you want to combine the two levels (like Stata defaults to), you can simply rename the columns:

```
In [13]: wide_single_level_column = wide.copy()

In [14]: wide_single_level_column.columns = [
...:     '{}_{}'.format(var, time)
...:     for var, time in wide_single_level_column.columns]

In [15]: wide_single_level_column
Out[15]:
```

	some_variable_1	some_variable_2
unit_id		
a	0	1
b	2	3
c	4	5
d	6	7

The `pivot` command can also be useful, but it's a bit more complicated than `stack` and `unstack` and is better to revisit `pivot` after you are comfortable working with DataFrame indexes and columns.

Stata	Python
<code>reshape <wide/long> <stubs>, i(<vars>)</code>	wide: <code>df.unstack(<level>)</code>
<code>j(<var>)</code>	long: <code>df.stack(<column_level>)</code>
	see also <code>df.pivot</code>

9 Econometrics

Stata	Python
<code>ttest <var>, by(<var>)</code>	<code>from scipy.stats import ttest_ind ttest_ind(<array1>, <array2>)</code>
<code>xi: i.<var></code>	<code>pd.get_dummies(df[<var>])</code>
<code>i.<var2>#c.<var1></code>	<code>pd.get_dummies(df[<var2>]).multiply(df[<var1>])</code>
<code>reg <yvar> <xvar> if <condition>, r</code>	<code>import econtools.metrics as mt results = mt.reg(df[<condition>], <yvar>, <xvar>, robust=True)</code>
<code>reg <yvar> <xvar> if <condition>, vce(cluster <clustervar>)</code>	<code>results = mt.reg(df[<condition>], <yvar>, <xvar>, cluster=<clustervar>)</code>
<code>areg <yvar> <xvar>, absorb(<fe_var>)</code>	<code>results = mt.reg(df, <yvar>, <xvar>, a_name= <fe_var>)</code>
<code>predict <newvar>, resid</code>	<code><newvar> = results.resid</code>
<code>predict <newvar>, xb</code>	<code><newvar> = results.yhat</code>
<code>_b[<var>], _se[<var>]</code>	<code>results.beta[<var>], results.se[<var>]</code>
<code>test <varlist></code>	<code>results.Ftest(<varlist>)</code>
<code>test <varlist>, equal</code>	<code>results.Ftest(<varlist>, equal=True)</code>
<code>lincom <var1> + <var2></code>	<code>econtools.metrics.f_test</code> with appropriate parameters.
<code>ivreg2</code>	<code>econtools.metrics.ivreg</code>
<code>outreg2</code>	<code>econtools.outreg</code>
<code>reghdfe</code>	None (hoping to add it to Econtools soon).

10 Plotting

Visualizations are best handled by the packages Matplotlib and Seaborn.

Stata	Python
<code>binscatter</code>	<code>econtools.binscatter</code>
<code>maptile</code>	No quick tool, but easy to do with Cartopy.
<code>coefplot</code>	<code>ax.scatter(results.beta.index, results.beta)</code> often works. Depends on context.

Stata	Python
<code>twoway scatter <var1> <var2></code>	<code>df.scatter(<var2>, <var1>)</code>
<code>twoway scatter <var1> <var2> if <condition></code>	<code>df[<condition>].scatter(<var2>, <var1>)</code>
<code>twoway <connected/line/area/bar/rarea></code>	As above, though <code>ax.plot(<var1>, <var2>)</code> is better. Like merge, it's a different paradigm, needs more explanation.

11 Other differences

11.1 Missing values

In Python, missing values are represented by a NumPy "not a number" object, `np.nan`. In Stata, missing (.) is larger than every number, so `10 < .` yields True. In Python, `np.nan` is never equal to anything. Any comparison involving `np.nan` is always False, even `np.nan == np.nan`.

To look for missing values in DataFrame columns, use any of the following.

- `df[<varname>].isnull()` returns a vector of True and False values for each row of `df[<varname>]`.
- `df[<varname>].notnull()` is the complement of `.isnull()`.
- The function `np.isnan(<arraylike>)` takes an array and returns True or False for each element of the array (a DataFrame is a special type of array).

Another important difference is that `np.nan` is a floating point data type, so any column of a DataFrame that contains missing numbers will be floats. If a column of integers gets changed so that even one row is `np.nan`, the whole column will be converted to floats.

11.2 Floating point equality

In Stata, decimal numbers are never equal to anything, e.g., `3.0 == 3` is False. This is not a problem in Python, the above equality check returns True.