

XVI Convegno Italiano degli Utenti di Stata - Firenze, 26-27 September,
2019

Calling Python Scripts in Stata: a Power-Law application

Antonio Zinilli
National Research Council (CNR)



Aim

Stata allows to call external routines, written in other software, to perform specific tasks within Stata.

This talk offers some insights on how to develop a Stata procedure embedding an external software routine (Python, in this case).

We provide a user-written Stata module python, written to allow users to run scripts by calling back Python programming language.



Motivation

The -python- package facilitates integrating Python with Stata by allowing automatic interprocess communication between the two software. Here we present a statistical application implemented in Python and called in Stata for discerning and quantifying power-law behavior in empirical data.

In a nutshell

- Call Python from Stata
- Output Python results within Stata
- Invoke Python interactively in do-files and ado-files
- Execute a Python script file(.py) directly through Stata



Motivation

Mix Stata code and Python code in a single do-file. We can take advantage of Stata's data management capabilities and packages implemented in Python (and not yet in Stata).

There are things we can do in Python that would be easier in Stata (e.g. graphs, descriptive analysis etc.).

Run multiple Python statements.



Case study: a Power-Law application

Data: citation count of the 1120 top-cited physicists from 1998-2014 (Source: Scopus)

Power Law function implements both the discrete and continuous maximum likelihood estimators for fitting the power-law distribution to data, along with the goodness-of-fit based approach to estimating the lower cut-off for the scaling region.

The probability $P(x)$ that the citations decays with the law:

$$P(x) \propto x^{-\alpha}$$

Where $P(x)$ is the probability to encounter value x and α as the scaling exponent ($\alpha > 0$)



5

Functional form by distribution

Distribution	$f(x)$
Log-normal	$\frac{1}{\sqrt{2\pi}^{\sigma_x}} e^{-(\ln x - \ln x_0)^2 / (2\sigma^2)}$
Exponential	$e^{-\lambda x}$
Stretched exponential (Weibull)	e^{-x^β}
Power-law + cut off	$x^{-\alpha} e^{-\lambda x}$



6

Intro to Python from Stata

Download the latest version of Python through Anaconda

It is the standard platform for Python data science

python search ← Search for Python installations on the current system

set python_exec "C:\Users\Zinilli\Anaconda3\python.exe", permanently

Set which version of Python to use

Python: ← Enter Python interactive environment

Issuing python (without a colon) will allow you to remain in the Python environment




7

Procedure for Power Law function

```
clear all
cd "C:\Users\Zinilli\Dropbox\Sug2019"
import excel "C:\Users\Zinilli\Dropbox\Sug2019\Citations.xlsx", sheet("Foglio1") firstrow
sum Citations

python:
from __future__ import division
import powerlaw
import numpy
import numpy as np
from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt
[!pip install package]
data=np.array(Data.get("Citations"))

results = powerlaw.Fit(data, fit_method='Likelihood', discrete=True, estimate_discrete=True)
```

Download the libraries

Additional package that is not part of the standard library

Function Power Law




8

Procedure for Power Law function

```

results.fixed_xmin
results.power_law.alpha
fit.power_law
fit.power_law.alpha
fit.power_law.parameter1
fit.lognormal.mu
fit.lognormal.parameter1
fit.exponential.parameter1
R,p=fit.loglikelihood_ratio('power_law', 'lognormal',normalized_ratio=True)
print ( R, p)
R,p=fit.loglikelihood_ratio('power_law', 'exponential',normalized_ratio=True)
print ( R, p)
R,p=fit.loglikelihood_ratio('power_law', 'truncated_power_law',nested=True)
print ( R, p)
R,p=fit.loglikelihood_ratio('power_law', 'stretched_exponential', nested=True)
print ( R, p)

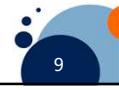
```

The search for the optimal xmin

α of the power law and its standard error, σ , are both calculated using the maximum likelihood method

Within the Fit object there are individual distribution objects for different distributions.
Each Distribution has the best fit parameters for that distribution.

R is the loglikelihood ratio between the two candidate distributions. This number will be positive if the data is more likely in the first distribution, and negative if the data is more likely in the second distribution.
The significance value for that direction is p.

9

Procedure for Power Law function

```

fig1 = fit.plot_ccdf(linewidth=3,label='Empirical Data')
fit.power_law.plot_ccdf(ax=fig1, color='black', linestyle='-.',label='Power Law', linewidth=0.8)
fit.lognormal.plot_ccdf(ax=fig1, color='g', linestyle='--',label='Lognormal')
fit.exponential.plot_ccdf(ax=fig1, color='y', linestyle=':',label='Exponential',linewidth=0.8)
fit.truncated_power_law.plot_ccdf(ax=fig1, color='m', linestyle=':', survival=True, label='Power Law + cut-off', linewidth=2)
fit.stretched_exponential.plot_ccdf(ax=fig1, color='b', linestyle='--',label='Stretched_exponential',linewidth=0.8)
fig1.set_ylabel(r'$p(X \geq x)$')
fig1.set_xlabel(r'Degree')
plt.xlabel("Strength Distribution -> Black = Empirical Data")
handles, labels = fig1.get_legend_handles_labels()
fig1.legend(handles, labels, loc=3, fontsize=9)
plt.title("Comparison of Power Law with Other Fits")
plt.savefig("powerlaw.png")

```

Closing Python environment

Via the library matplotlib. Function to create a single plot.




10

Screenshot

```

1 clear all
2 * Import Data
3 import excel "C:\Users\zilizili\Desktop\bipolar.xlsx", sheet("Foglio1") firstrow
4 python
5 future _import division
6 import powerlaw
7 import numpy as np
8 import matplotlib.pyplot as plt
9 import matplotlib3d import Axes3D
10 import matplotlib.pyplot as plt
11 numpy.seterr(divide='ignore')
12 data=np.array(Data.get("Citations"))
13 data=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True)
14 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
15 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
16 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
17 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
18 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
19 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
20 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
21 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
22 results=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
23 print(results.xmin)
24 print(results.xmax)
25 print(results.alpha)
26 print(results.powerlaw_alpha)
27 results = powerlaw.Fit(data,discrete=True,estimate_discrete=True,discrete_approximation='round')
28 print(results.fixed_xmin)
29 print(results.fixed_xmax)
30 print(results.powerlaw_alpha)
31 print(results.powerlaw_alpha)
32 print(results.powerlaw_alpha)
33 print(results.powerlaw_alpha)
34 print(results.powerlaw_alpha)
35 print(results.powerlaw_alpha)
36 print(results.powerlaw_alpha)
37 print(results.powerlaw_alpha)
38 print(results.powerlaw_alpha)
39 print(results.powerlaw_alpha)
40 print(results.powerlaw_alpha)
41 print(results.powerlaw_alpha)
42 print(results.powerlaw_alpha)
43 print(results.powerlaw_alpha)
44 print(results.powerlaw_alpha)
45 print(results.powerlaw_alpha)
46 print(results.powerlaw_alpha)
47 print(results.powerlaw_alpha)
48 print(results.powerlaw_alpha)
49 print(results.powerlaw_alpha)
50 print(results.powerlaw_alpha)
51 print(results.powerlaw_alpha)
52 print(results.powerlaw_alpha)
53 print(results.powerlaw_alpha)
54 print(results.powerlaw_alpha)
55 print(results.powerlaw_alpha)
56 print(results.powerlaw_alpha)
57 print(results.powerlaw_alpha)
58 print(results.powerlaw_alpha)
59 print(results.powerlaw_alpha)
60 print(results.powerlaw_alpha)
61 print(results.powerlaw_alpha)
62 print(results.powerlaw_alpha)
63 print(results.powerlaw_alpha)
64 print(results.powerlaw_alpha)
65 end

```

RCFES 11

Python output visible as Stata output

```

Statistics User Window Help
File Edit View Language Project Tools
Do-file Editor - Power_law.do
Power_law.do*
1 >>> import powerlaw
2 >>> import numpy
3 >>> import numpy as np
4 >>> from mpl_toolkits.mplot3d import Axes3D
5 >>> import matplotlib.pyplot as plt
6 >>> numpy.seterr(divide='ignore', invalid='ignore')
7 >>> data=np.array(Data.get("Citations"))
8 >>> data=powerlaw.Fit(data,fit_method='Likelihood',discrete=True,estimate_discrete=True,discrete_approximation='round')
9 Calculating best minimal value for power law fit
10 >>> results=xmin
11 3712.0
12 >>> results.fixed_xmin
13 False
14 >>> results.power_law.alpha
15 4.0762955240709055
16 >>> results.power_law.D
17 0.0313679699652151
18 >>> print(results.xmin)
19 3712.0
20 >>> nprint(results.fixed_xmin)

```

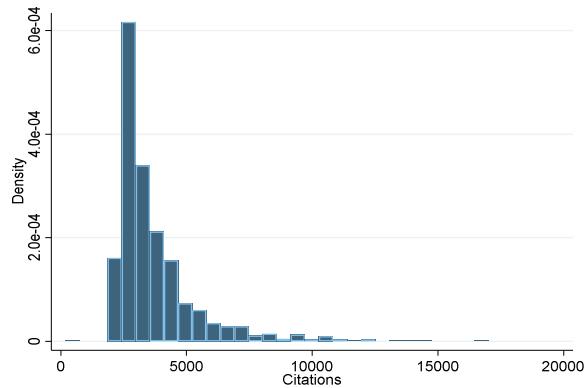
RCFES 12

Descriptive analysis

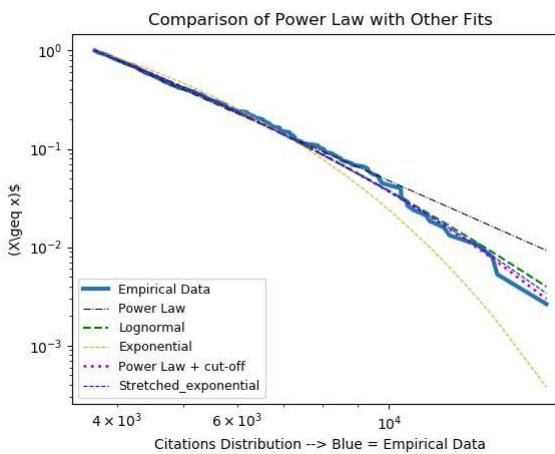
. sum Citations

Variable	Obs	Mean	Std. Dev.	Min	Max
Citations	1,119	3705.532	1690.558	178	16994

. twoway (histogram Citations)



Complementary cumulative distribution functions of example data and fitted distributions

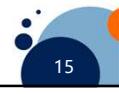


Power-law and alternative distributions

alpha= 4
xmin=3712

$$P(x) \propto x^{-\alpha}$$

	Power law	Log-normal		Exponential		Stretched exponential		Power-law + cut off	
Citations	p*	LR	p*	LR	p*	LR	p*	LR	p*
	0,03	-1.32	0,18	0,72	0,47	-2,23	0,03	-2,55	0,02



Conclusions

Summarizing, we cannot conclude on the exact nature of the far-tail of distributions of citations.

Power Law with cut off and Stretched exponential seems to be both good, we need further analysis to investigate the best.

