

Linear models

Friday, September 20

Content

- [1. The statsmodel package](#)
- [2. OLS](#)
- [3. Generalised linear model](#)
- [4. Two-stage least squares](#)

1. The statsmodels package

statsmodels is a Python module that provides classes and functions for the estimation of many different statistical models, as well as for conducting statistical tests, and statistical data exploration.

The online documentation is hosted at [statsmodels.org](https://www.statsmodels.org/stable/index.html) (<https://www.statsmodels.org/stable/index.html>)

It covered:

- Linear Regression
- Generalized Linear Models
- Generalized Estimating Equations
- Generalized Additive Models (GAM)
- Robust Linear Models
- Linear Mixed Effects Models
- Regression with Discrete Dependent Variable
- Generalized Linear Mixed Effects Models
- ANOVA
- Time Series analysis `tsa`
- Time Series Analysis by State Space Methods `statespace`
- Vector Autoregressions `tsa.vector_ar`
- Methods for Survival and Duration Analysis
- Statistics `stats`
- Nonparametric Methods `nonparametric`
- Generalized Method of Moments `gmm`
- Contingency tables
- Multiple Imputation with Chained Equations
- Multivariate Statistics `multivariate`
- Empirical Likelihood `emplike`
- Other Models `miscmodels`
- Distributions
- Graphics
- Input-Output `iolib`
- Tools
- The Datasets Package
- Sandbox
- Working with Large Data Sets
- Optimization

`statsmodels` works smoothly with the `pandas` in a way that `DataFrame` is the dataset form it supports by default.

Anaconda has installed `statsmodels` module by default. Before using the functions and classes inside, we need to import the `statsmodels.api` and `statsmodels.formula.api`.

In [1]:

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

import numpy as np
```

The output of `statsmodels` is similar to the output of functions in `R`. We start with the most widely used and elementary statistical methods : ordinary least square.

2. OLS

2.1. How to fit a dataset and see the result

We use the dataset `Guerry` provided by `statsmodel` which studied the determinants of the number of lottery sold.

In [2]:

```
import pandas as pd
# Load data
dat = sm.datasets.get_rdataset("Guerry", "HistData").data
# List of the variables
dat.head(5)
```

Out[2]:

	dept	Region	Department	Crime_pers	Crime_prop	Literacy	Donations	Infants	Suicides
0	1	E	Ain	28870	15890	37	5098	33120	35036
1	2	N	Aisne	26226	5521	51	8901	14572	12837
2	3	C	Allier	26747	7925	13	10973	17044	11412
3	4	E	Basses-Alpes	12935	7289	46	2733	23018	14238
4	5	E	Hautes-Alpes	17488	8174	69	6962	23076	16177

5 rows × 23 columns

More specifically, we studied the relationship between lottery and the literacy and population (in the log scale).

In [3]:

```
# Fit regression model (using the natural log of one of the regressors)
model = smf.ols('Lottery ~ Literacy + np.log(Pop1831)', data=dat)
results = model.fit()
```

To see the results, we need an additional step:

In [4]:

```
# Inspect the results
print(results.summary())
```

```

                                OLS Regression Results
=====
====
Dep. Variable:                  Lottery    R-squared:
0.348
Model:                          OLS      Adj. R-squared:
0.333
Method:                        Least Squares    F-statistic:                2
2.20
Date:                          Thu, 12 Sep 2019    Prob (F-statistic):            1.90
e-08
Time:                          22:19:49    Log-Likelihood:                -37
9.82
No. Observations:                86    AIC:                            7
65.6
Df Residuals:                    83    BIC:                            7
73.0
Df Model:                        2
Covariance Type:                nonrobust
=====
=====
                                coef      std err          t      P>|t|      [0.025
0.975]
-----
-----
Intercept                246.4341      35.233        6.995      0.000      176.358
316.510
Literacy                 -0.4889       0.128       -3.832      0.000      -0.743
-0.235
np.log(Pop1831)         -31.3114       5.977       -5.239      0.000     -43.199
-19.424
=====
====
Omnibus:                  3.713    Durbin-Watson:
2.019
Prob(Omnibus):            0.156    Jarque-Bera (JB):
3.394
Skew:                     -0.487    Prob(JB):
0.183
Kurtosis:                 3.003    Cond. No.
702.
=====
=====

```

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

2.2. When the dataset is not in DataFrame

The dataset above is provided by `statsmodels` package hence in the form it supports. However, in many situations, the dataset is not constructed yet. In this case, we can use `numpy` arrays.

In [5]:

```
import numpy as np

import statsmodels.api as sm

# Generate artificial data (2 regressors + constant)
nobs = 100

X = np.random.random((nobs, 2))

X = sm.add_constant(X)

beta = [1, .1, .5]

e = np.random.random(nobs)

y = np.dot(X, beta) + e

# Fit regression model
results = sm.OLS(y, X).fit()

# Inspect the results
print(results.summary())
```

OLS Regression Results

```

=====
====
Dep. Variable:          y    R-squared:
0.164
Model:                OLS    Adj. R-squared:
0.146
Method:              Least Squares    F-statistic:
9.495
Date:                Thu, 12 Sep 2019    Prob (F-statistic):        0.00
0171
Time:                22:19:50    Log-Likelihood:            -1
0.743
No. Observations:        100    AIC:                2
7.49
Df Residuals:            97    BIC:                3
5.30
Df Model:                2
Covariance Type:        nonrobust
=====
====
               coef    std err          t      P>|t|      [0.025    0.
975]
-----
----
const          1.6012     0.076    21.023     0.000     1.450
1.752
x1            -0.0612     0.099     -0.619     0.538    -0.257
0.135
x2             0.4040     0.097     4.183     0.000     0.212
0.596
=====
====
Omnibus:            11.587    Durbin-Watson:
1.954
Prob(Omnibus):      0.003    Jarque-Bera (JB):
4.146
Skew:               0.156    Prob(JB):
0.126
Kurtosis:           2.053    Cond. No.
5.67
=====
====

```

Warnings:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

Of course, we can create a dataset and make it supported by `statsmodels`. Details can be found here: [adding a dataset \(https://www.statsmodels.org/stable/dev/dataset_notes.html?highlight=statsmodels%20datasets#adding-a-dataset-an-example\)](https://www.statsmodels.org/stable/dev/dataset_notes.html?highlight=statsmodels%20datasets#adding-a-dataset-an-example)

2.3. Wald's test

Besides the fitting, `statsmodels` also supports many statistical testing methods. Here, we show how to use *Wald's test* in `statsmodels`.

Again, we consider the dataset `Guerry`.

We want to analyse the effect of *Wealth* and *Literacy* on the `_Crimepers` and test:

whether the coefficients of *Wealth* and *Literacy* are the same.

In [6]:

```
formula = 'Crime_pers ~ Wealth + Literacy'
results = smf.ols(formula, dat).fit()
hypotheses = '(Wealth = Literacy)'
f_test = results.f_test(hypotheses)
print(f_test)
```

```
<F test: F=array([[0.03467668]]), p=0.8527291641569565, df_denom=83, df_nu
m=1>
```

3. Generalised linear model

Generalized linear models in `statsmodels` currently supports estimation using the one-parameter exponential families.

What is it?

The statistical model for each observation i is assumed to be

$$Y_i \sim F_{EDM}(\cdot | \theta, \phi, w_i) \text{ and } \mu_i = E[Y_i | x_i] = g^{-1}(x_i' \beta).$$

where g is the link function and $F_{EDM}(\cdot | \theta, \phi, w)$ is a distribution of the family of exponential dispersion models (EDM) with natural parameter θ , scale parameter ϕ and weight w . Its density is given by

$$f_{EDM}(y | \theta, \phi, w) = c(y, \phi, w) \exp\left(\frac{y\theta - b(\theta)}{\phi} w\right).$$

It follows that $\mu = b'(\theta)$ and $Var[Y | x] = \frac{\phi}{w} b''(\theta)$. The inverse of the first equation gives the natural parameter as a function of the expected value $\theta(\mu)$ such that

$$Var[Y_i | x_i] = \frac{\phi}{w_i} v(\mu_i)$$

with $v(\mu) = b''(\theta(\mu))$. Therefore it is said that a GLM is determined by link function g and variance function $v(\mu)$ alone (and x of course).

Note that while ϕ is the same for every observation y_i and therefore does not influence the estimation of β , the weights w_i might be different for every y_i such that the estimation of β depends on them.

Examples: binomial response $B(n, p)$

$$\begin{aligned} \mu &= E[Y | x] = np \\ v(\mu) &= \mu - \frac{\mu^2}{n} \\ \theta(\mu) &= \log \frac{p}{1-p} \\ b(\theta) &= n \log(1 + e^\theta) \end{aligned}$$

Real data

Here we use the Star98 dataset which was taken with permission from Jeff Gill (2000) *Generalized linear models: A unified approach*

In [7]:

```
import statsmodels.api as sm
import statsmodels.formula.api as smf

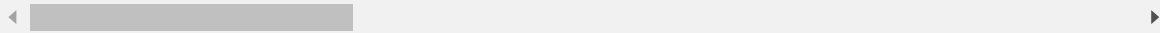
import pandas as pd

star98 = sm.datasets.star98.load(as_pandas=True)
star98.data.head(5)
```

Out[7]:

	NABOVE	NBELOW	LOWINC	PERASIAN	PERBLACK	PERHISP	PERMINTE	AVYRSEXP
0	452.0	355.0	34.39730	23.299300	14.235280	11.411120	15.91837	14.70646
1	144.0	40.0	17.36507	29.328380	8.234897	9.314884	13.63636	16.08324
2	337.0	234.0	32.64324	9.226386	42.406310	13.543720	28.83436	14.59559
3	395.0	178.0	11.90953	13.883090	3.796973	11.443110	11.11111	14.38939
4	8.0	57.0	36.88889	12.187500	76.875000	7.604167	43.58974	13.90568

5 rows × 22 columns



In [8]:

```
print(sm.datasets.star98.NOTE)
```

```
::
```

Number of Observations - 303 (counties in California).

Number of Variables - 13 and 8 interaction terms.

Definition of variables names::

the
the
y the
umber

- NABOVE - Total number of students above the national median for
math section.
- NBELOW - Total number of students below the national median for
math section.
- LOWINC - Percentage of low income students
- PERASIAN - Percentage of Asian student
- PERBLACK - Percentage of black students
- PERHISP - Percentage of Hispanic students
- PERMINTE - Percentage of minority teachers
- AVYRSEXP - Sum of teachers' years in educational service divided b
number of teachers.
- AVSALK - Total salary budget including benefits divided by the n
of full-time teachers (in thousands)
- PERSPENK - Per-pupil spending (in thousands)
- PTRATIO - Pupil-teacher ratio.
- PCTAF - Percentage of students taking UC/CSU prep courses
- PCTCHRT - Percentage of charter schools
- PCTYRRND - Percentage of year-round schools

The below variables are interaction terms of the variables defined above.

- PERMINTE_AVYRSEXP
- PEMINTE_AVSAL
- AVYRSEXP_AVSAL
- PERSPEN_PTRATIO
- PERSPEN_PCTAF
- Ptratio_PCTAF
- PERMINTE_AVYRSEXP_AVSAL
- PERSPEN_PTRATIO_PCTAF

Now, we use generalized linear model to analyze the number of students above the national median for the math section

In [9]:

```
data = sm.datasets.star98.load(as_pandas=False)
data.exog = sm.add_constant(data.exog, prepend=False)
glm_binom = sm.GLM(data.endog, data.exog, family=sm.families.Binomial())
res = glm_binom.fit()
print(res.summary())
```

Generalized Linear Model Regression Results

```

=====
====
Dep. Variable:          ['y1', 'y2']    No. Observations:
303
Model:                  GLM             Df Residuals:
282
Model Family:          Binomial         Df Model:
20
Link Function:         logit            Scale:                1.
0000
Method:                IRLS             Log-Likelihood:       -29
98.6
Date:                  Thu, 12 Sep 2019   Deviance:              40
78.8
Time:                  22:19:56          Pearson chi2:           4.05
e+03
No. Iterations:        5
Covariance Type:       nonrobust
=====
=====

```

```

=====
====
              coef      std err          z      P>|z|      [0.025      0.
975]
-----
----
x1          -0.0168      0.000     -38.749      0.000     -0.018      -
0.016
x2           0.0099      0.001      16.505      0.000      0.009
0.011
x3          -0.0187      0.001     -25.182      0.000     -0.020      -
0.017
x4          -0.0142      0.000     -32.818      0.000     -0.015      -
0.013
x5           0.2545      0.030       8.498      0.000      0.196
0.313
x6           0.2407      0.057       4.212      0.000      0.129
0.353
x7           0.0804      0.014       5.775      0.000      0.053
0.108
x8          -1.9522      0.317      -6.162      0.000     -2.573      -
1.331
x9          -0.3341      0.061      -5.453      0.000     -0.454      -
0.214
x10         -0.1690      0.033      -5.169      0.000     -0.233      -
0.105
x11          0.0049      0.001       3.921      0.000      0.002
0.007
x12         -0.0036      0.000     -15.878      0.000     -0.004      -
0.003
x13         -0.0141      0.002      -7.391      0.000     -0.018      -
0.010
x14         -0.0040      0.000      -8.450      0.000     -0.005      -
0.003
x15         -0.0039      0.001      -4.059      0.000     -0.006      -
0.002
x16          0.0917      0.015       6.321      0.000      0.063
0.120
x17          0.0490      0.007       6.574      0.000      0.034
0.064
x18          0.0080      0.001       5.362      0.000      0.005
0.011

```

x19	0.0002	2.99e-05	7.428	0.000	0.000	
0.000						
x20	-0.0022	0.000	-6.445	0.000	-0.003	-
0.002						
const	2.9589	1.547	1.913	0.056	-0.073	
5.990						
=====						
=====						
=====						
=====						

Also, we can see the fit plot of the glm model

In [10]:

```
nobs = res.nobs
y = data.endog[:,0]/data.endog.sum(1)
yhat = res.mu

from statsmodels.graphics.api import abline_plot

from matplotlib import pyplot as plt

fig, ax = plt.subplots()
ax.scatter(yhat, y)
line_fit = sm.OLS(y, sm.add_constant(yhat, prepend=True)).fit()
abline_plot(model_results=line_fit, ax=ax)

ax.set_title('Model Fit Plot')
ax.set_ylabel('Observed values')
ax.set_xlabel('Fitted values');
```

4. Two-stage least squares

4.1. Endogeneity

Endogeneity issues are at the central of the quantitative research in the social science. That is to say, when we use the linear regression, the dependent variable might actually affect the explanatory variable. And once this happens, the estimates from the OLS could be largely biased.

For example, there is a two-way relationship between the institutions and the economic outcomes:

- better institutions will output labor force of higher quality which boost the economic development
- richer countries/cities can afford better institutions

To eliminate such endogeneity, two-stage least square method is one tool used by many social scientists. The idea is to find an *instrument variable* that is

- correlated with the explanatory variable
- not correlated with the dependent variable

4.2. Real data: Acemoglu et al. (2001)

As an example, we will use the data set from Daron Acemoglu, Simon Johnson, and James A Robinson. *The colonial origins of comparative development: an empirical investigation*. The American Economic Review, 91(5):1369–1401, 2001.

In this paper, Acemoglu et al. (2001) want to study the effect of the institution quality on the economic outcomes.

The data set could be downloaded from [Quant Econ \(https://lectures.quantecon.org/\)](https://lectures.quantecon.org/)

In [11]:

```
import pandas as pd

# Import and select the data
df4 = pd.read_stata('https://github.com/QuantEcon/QuantEcon.lectures.code/raw/master/ols/maketable4.dta')
df4 = df4[df4['baseco'] == 1]

df4.head(5)
```

Out[11]:

	shortnam	africa	lat_abst	rich4	avexpr	logpgp95	logem4	asia	loghjypl	basecc
1	AGO	1.0	0.136667	0.0	5.363636	7.770645	5.634789	0.0	-3.411248	1.0
3	ARG	0.0	0.377778	0.0	6.386364	9.133459	4.232656	0.0	-0.872274	1.0
5	AUS	0.0	0.300000	1.0	9.318182	9.897972	2.145931	0.0	-0.170788	1.0
11	BFA	1.0	0.144444	0.0	4.454545	6.845880	5.634789	0.0	-3.540459	1.0
12	BGD	0.0	0.266667	0.0	5.136364	6.877296	4.268438	1.0	-2.063568	1.0

Acemoglu et al. (2001) use:

- economic outcome: *logpgp95* , log GDP per capita in 1995, adjusted for exchange rates
- institution quality: *avexpr* , an index of protection against expropriation on average over 1985-95
- instrument variable: *logem4* , settler mortality rates

In [12]:

```
import statsmodels.sandbox.regression.gmm as gmm

model = gmm.IV2SLS(endog=df4['logpgp95'], exog=df4['avexpr'], instrument=df4['logem4'])
result = model.fit()
print(result.summary())
```

IV2SLS Regression Results

```
=====
====
Dep. Variable:          logpgp95    R-squared:
0.976
Model:                  IV2SLS      Adj. R-squared:
0.975
Method:                 Two Stage    F-statistic:
nan
                                Least Squares    Prob (F-statistic):
nan
Date:                   Thu, 12 Sep 2019
Time:                   22:20:00
No. Observations:      64
Df Residuals:          63
Df Model:               1
=====
====
                coef    std err          t      P>|t|      [0.025    0.
975]
-----
----
avexpr          1.2468      0.026     47.531      0.000      1.194
1.299
=====
====
Omnibus:              0.340    Durbin-Watson:
2.052
Prob(Omnibus):        0.844    Jarque-Bera (JB):
0.474
Skew:                 -0.152    Prob(JB):
0.789
Kurtosis:             2.707    Cond. No.
1.00
=====
====
```

```
C:\Users\dyevre\AppData\Local\Continuum\anaconda3\lib\site-packages\scipy
\stats\_distn_infrastructure.py:877: RuntimeWarning: invalid value encount
ered in greater
    return (self.a < x) & (x < self.b)
C:\Users\dyevre\AppData\Local\Continuum\anaconda3\lib\site-packages\scipy
\stats\_distn_infrastructure.py:877: RuntimeWarning: invalid value encount
ered in less
    return (self.a < x) & (x < self.b)
C:\Users\dyevre\AppData\Local\Continuum\anaconda3\lib\site-packages\scipy
\stats\_distn_infrastructure.py:1831: RuntimeWarning: invalid value encoun
tered in less_equal
    cond2 = cond0 & (x <= self.a)
```