

Python for Finance

Regressions, Interpolation & Optimisation

Andras Niedermayer



Outline

- ① Regressions in pandas
- ② Function approximation
 - Regression
 - Interpolation
- ③ Convex optimization

Pivot tables

We want to create a table with opening prices, using index names as columns.

We can use the pivot function in Pandas.

```
1 PivotOpen=IndicesA.pivot( index = 'Date',  
2     columns='Index', values='Open')
```

Running the pivot function without the *values* option creates a collection of tables. To select the 'Open':

```
1 PivotTable=IndicesA.pivot(index='Date',  
2     columns='Index')  
3 PivotOpen=PivotTable['Open']
```

Stacking and unstacking

First, let us index the data by date and index name:

```
1 IxNew.set_index(['Date', 'Index'], inplace=True)
```

To collapse the column of a database to a single (data) series:

```
1 IxNewStack=IxNew.stack()
```

To restore the indices as columns (sometimes useful):

```
1 IxNewStack=IxNew.stack().reset_index()
```

To restore the original database:

```
1 IxNew=IxNewStack.unstack()
```

OLS regression (statsmodels.api)

Suppose we want to regress *index returns* (just computed) on *index daily volatility* (high-low range).

$$\text{Return}_{it} = \alpha + \beta \frac{\text{High}_{it}}{\text{Low}_{it}} + \varepsilon_{it}$$

The simplest OLS model reads:

```
1 model = sm.OLS(Y, X)
2 results = model.fit()
3 print(results.summary())
```

We can access the results as:

- ① Coefficient estimates: `results.params`
- ② Estimator covariance matrix: `results.cov_HCO`
- ③ p-values: `results.pvalues`
- ④ R-squared: `results.rsquared`

OLS regression output

Dep. Variable:	Return	R-squared:	0.018
Model:	OLS	Adj. R-squared:	0.018
Method:	Least Squares	F-statistic:	115.6
Date:	Tue, 13 Feb 2018	Prob (F-statistic):	1.00e-26
Time:	23:59:34	Log-Likelihood:	19253.
No. Observations:	6345	AIC:	-3.850e+04
Df Residuals:	6343	BIC:	-3.849e+04
Df Model:	1		

	coef	std err	t	P> t	[0.025	0.975]
const	0.1841	0.017	10.765	0.000	0.151	0.218
0	-0.1813	0.017	-10.751	0.000	-0.214	-0.148
Omnibus:	680.608		Durbin-Watson:	1.987		
Prob(Omnibus):	0.000		Jarque-Bera (JB):	3744.882		
Skew:	0.367		Prob(JB):	0.00		
Kurtosis:	6.691		Cond. No.	234.		

Useful regression output

To see all regression data and not just the summary, type:
`model .+<Tab>`

Applications

Starting from the IxNew data:

- 1 On how many days was the return larger for CAC40 than for DAX?
- 2 Create a Series object indexed by Date that contain the name of the index with the highest return.

Hint: Use the `idxmax` method:

<http://pandas.pydata.org/pandas-docs/stable/generated/pandas.DataFrame.idxmax.html>

Applications - (One) Solution

1.

```
1 PivotReturn = IxNew.reset_index().pivot('Date',
2     'Index', 'Return')
3 days_a = (PivotReturn['CAC40']
4     > PivotReturn['DAX']).sum()
```

2.

```
1 PivotReturn.apply(lambda x: x.idxmax(), axis=1)
2 PivotReturn.idxmax(axis=1)
```

Outline

- ① Regressions in pandas
- ② **Function approximation**
 - Regression
 - Interpolation
- ③ Convex optimization

Motivation

- ① Most of the times in finance, we do not know the DGP (Data Generating Process).
- ② Many applications in finance involve “reverse engineering” patterns from data.
- ③ This is useful, for example to make predictions about the future dynamics of financial variables.
- ④ Two main techniques:
 - ① Regression
 - ② Interpolation

First, define a function (the DGP)...

We specifically choose a non-polynomial function (more difficult).

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 def f(x):
5     return np.sin(x)+0.5*x
```

Next, generate data from the DGP

We generate 50 data-points from the DGP: $(x, f(x))$.

- Function `x=np.linspace(a, b, N)` returns an array of N numbers, equally spaced, from a to b .
- What does $f(x)$ return?

```
1 x=np.linspace(-2*np.pi, 2*np.pi, 50)
2
3 plt.plot(x, f(x), 'b')
4 plt.grid()
5 plt.xlabel('x', fontsize=18)
6 plt.ylabel('y', fontsize=18)
7 plt.show()
```

Regression

Theoretical framework:

- 1 You are given N points in a 2-D (can be 3-D, 4-D...) space: (x_j, y_j) .
- 2 You choose K (base) functions of x_j , i.e., $b_i(x_j)$, such that you believe y_j can be written as a linear combination of these functions.
- 3 You select coefficients of said linear combinations, α_i by minimizing the squared difference from the actual data.

$$\min_{\alpha_i} \frac{1}{N} \sum_{j=1}^N \left(y_j - \sum_{i=1}^K \alpha_i b_i(x_j) \right)^2 \quad (1)$$

Polynomial regression

A simple case is to approximate y_j as a polynomial function of x_j .

That is, choose: $b_0 = 1$, $b_1 = x$, $b_2 = x^2$, ..., $b_k = x^k$.

Easy to implement in Python with `polyfit` (polynomial fit):

- 1 First, get the coefficient list using `polyfit`.
- 2 Next, get the fitted values from the coefficient list using `polyval`.

```
1 reg=np.polyfit(x, f(x), deg=k)
2 y_fit=np.polyval(reg, x)
```

What happens if we vary the polynomial degree?

Polynomial regression

Application: Beyond polynomials

- The mean squared error of our fit is not zero....rather 1.77×10^{-3} .
 - Not surprising, since the original function was not a polynomial.
 - How can we approximate it using other base functions, i.e., trigonometric?
- ① Say we know (prior theoretical work) our function is a combination of a second order polynomial and sin/cos functions.
 - ② Let us define a matrix with values for $1, x, x^2, \sin(x), \cos(x)$

Application: Formalization of the problem

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_N \end{pmatrix} = \underbrace{\begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \dots \\ \alpha_N \end{pmatrix}}_{\text{Coefficients}} \underbrace{\begin{pmatrix} 1 & x_1 & x_1^2 & \sin(x_1) & \cos(x_1) \\ 1 & x_2 & x_2^2 & \sin(x_2) & \cos(x_2) \\ 1 & x_3 & x_3^2 & \sin(x_3) & \cos(x_3) \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_N & x_N^2 & \sin(x_N) & \cos(x_N) \end{pmatrix}}_{\text{Matrix M}} + \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ \dots \\ u_N \end{pmatrix}$$

Application: Solving the problem in Python (1/2)

- Initialize the matrix M:

```
1 matrix=np.zeros((len(x),5))
```

- Fill in each column with a variable:

```
1 matrix[:,0]=1
2 matrix[:,1]=x
3 matrix[:,2]=x**2
4 matrix[:,3]=np.sin(x)
5 matrix[:,4]=np.cos(x)
```

Application: Solving the problem in Python (2/2)

We use `numpy.linalg.lstsq` to minimise the sum of squared residuals.

Least-square coefficients are given by:

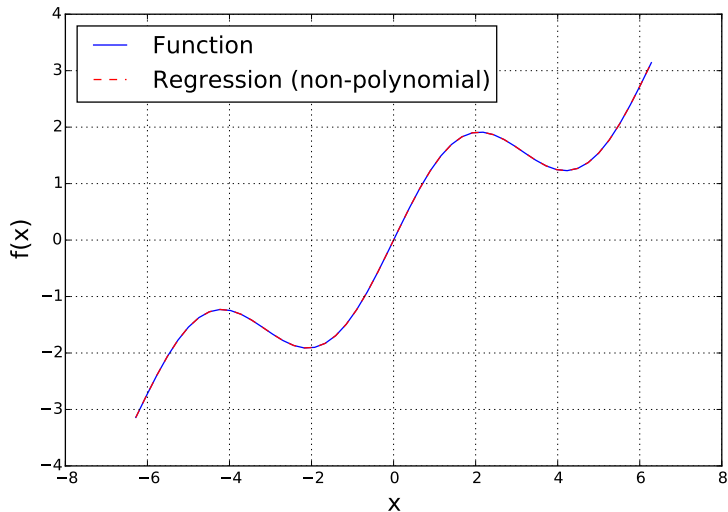
```
1 reg=np.linalg.lstsq(M, f(x))[0]
```

The fitted-values are computed as a dot-product between the coefficients vector (`reg`) and the matrix M :

```
1 y_fit2=np.dot(reg,M.T)
2 # we need to transpose the matrix
```

- ① What are the coefficients in `reg`?
- ② What is the MSE?

Application: Output



General idea

- ① With *regression*, one tries to identify a **unique** function $g(x)$ that is as close as possible to the “true”, unknown function $f(x)$, i.e.,

$$\min \sum (g(x) - f(x))^2$$

- ② With *interpolation*, one fits more (generally polynomial) functions, one between **each pair of consecutive points**.
- The fit is perfect, i.e., $\forall i, g_i(x_i) = f(x_i)$.
 - The function is not unique, which is mathematically involved.
 - The function is constrained to be continuous, $g_i(x_i) = g_{i+1}(x_i)$.
 - Some additional constraint is needed, i.e., second derivatives are continuous.
- ③ One needs ordered data in interpolation (unlike regression).
- ④ Procedure takes more time and is less parsimonious (more coefficients in the end) – but generally more accurate.

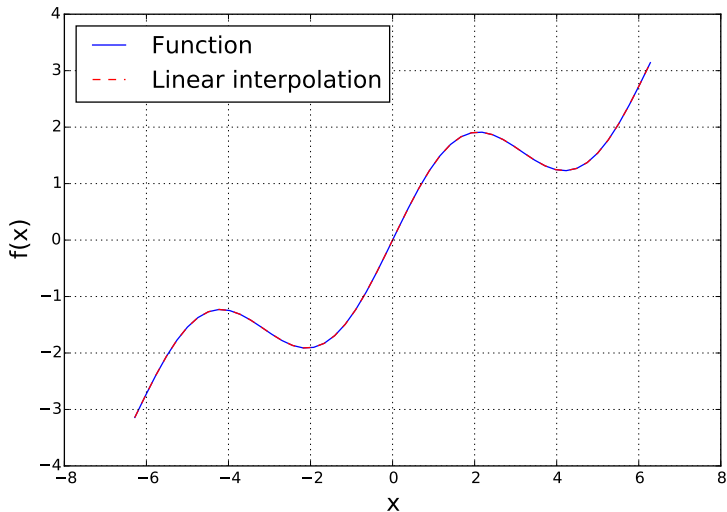
Implementation

The interpolation package is in the Scientific Python library (`scipy`). The parameter k defines the degree of the polynomial ($k = 1$ is a linear spline, $k = 3$ a cubic spline...)

```
1 import scipy.interpolate as spi
2 interp=spi.splrep(x,f(x), k=1)
3 y_interp=spi.splev(x,interp)
```

- ① What type of object is `interp` relative to `reg`? Why?
- ② How good is the linear interpolation?

Interpolation output



Outline

- ① Regressions in pandas
- ② Function approximation
 - Regression
 - Interpolation
- ③ Convex optimization

Main idea

We want to minimize a function $f(x_1, x_2, x_3, \dots, x_n)$:

$$\min_{x_i} f(x_1, x_2, x_3, \dots, x_n) \quad (2)$$

All local extrema satisfy

$$\frac{\partial f}{\partial x_i} = 0, \forall i \in \{1, 2, \dots, n\}. \quad (3)$$

The global minimum/maximum (if it exists and/or is unique) is either one of the local extrema or one of the domain end-points (see whiteboard).

More?

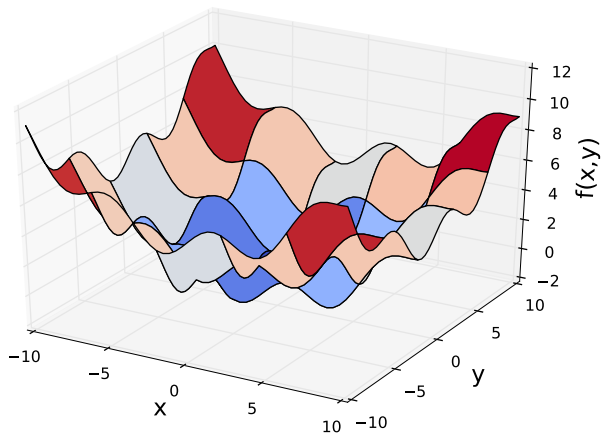
The Weierstrass (extreme value) theorem guarantees the existence of a maximum and minimum on closed and bounded intervals.

A two dimensional function

First, we define a function to minimize

```
1 def fm((x,y)):  
2     return np.sin(x)+1/20.0*x**2  
3     +np.sin(y)+1/20.0*y**2
```

A three dimensional graphic



Brute force optimization (the “caveman” approach)

```
1 import scipy.optimize as spo
```

Define a range and step to search for minimum:

```
1 search_area=(-10,10.01,5)
```

Change the function to print all iterations and output:

```
1 def fm((x,y)):  
2     z=np.sin(x)+1/20.0*x**2  
3     +np.sin(y)+1/20.0*y**2  
4     print '%8.4f□%8.4f□%8.4f' %(x,y,z)  
5     return z
```

Run the function brute (force) to find the minimum:

```
1 min_1=spo.brute(fm, (search_area,search_area),  
2     finish=None)
```

Brute force optimization (the “caveman” approach)

- ① What is the minimum found by this method?
- ② How can we improve the accuracy? What is the drawback?

The brute force method, while limited, can serve to provide starting values for more sophisticated algorithms.

One such function, working with numerical gradients, is `fmin`.

Optimization with `fmin`

General structure:

```
1 [xopt, fopt]=spo.fmin(function, start_values,  
2   xtol=, ftol=, maxiter=, maxfun=,)
```

- ① `xtol` : Relative error in argument acceptable for convergence.
- ② `ftol`: Relative error in function acceptable for convergence.
- ③ `maxiter` : Maximum number of iterations to perform.
- ④ `maxfun` : Maximum number of function evaluations to make.

We can use the global optimization results as starting values:

```
1 min_2=spo.fmin(fm, min_1, xtol=0.001, ftol=0.001)
```

Caveats

- Local optimization routines can get stuck in local extrema...
- ... or they may never converge.
- It is a good idea to perform a global optimization first to pinpoint the neighborhood of global minimum.
- What happens if we start `fmin` with $(2, 2)$ as starting values?

Constrained optimization

Most of the time, we look for optimal values of a function **under a set of constraints**.

Problem

There are two securities, A and B: Both cost 10 today. Tomorrow there are two equally likely states of the world: g or b . In state g , $A = 15$ and $B = 5$. In state b , $A = 5$ and $B = 12$. Assume an investor has 100 units of cash today and utility $u(w) = \sqrt{w}$. What is his optimal investment?

Application: Formalization of the problem

$$\max_{a,b} \mathbb{E}u(w_1) = \max_{a,b} \frac{1}{2}\sqrt{15a + 5b} + \frac{1}{2}\sqrt{5a + 12b}, \quad (4)$$

subject to:

$$10a + 10b \leq 100. \quad (5)$$

Application: Python implementation (Solution)

First, define the function. Note: we want to **maximize** expected utility!

```
1 def ExpU((s, b)):  
2     return -(0.5*np.sqrt(s*15+b*5)+  
3     0.5*np.sqrt(s*5+b*12))
```

Second, define the constraint as a dict variable and an implicit function. Inequality sign is always implicitly " ≥ 0 ".

```
1 cons = ({ 'type': 'ineq', 'fun':  
2     lambda (s, b): 100 - s*10 - b*10})
```

Third, choose starting values:

```
1 startval = [5, 5]
```

Fourth, run the minimize function from the optimization package:

```
1 result = spo.minimize(ExpU, startval,  
~     method='SLSQP', constraints=cons)
```

Notes

- method stands for optimization algorithm. SLSQP (Sequential Least Squares Programming) allows one to introduce constraints.
- One can specify Jacobian (`jac`) or Hessian matrix (`hess`) directly.
- In addition, bounds for the argument can be specified by `bounds`.

Output methods:

- ① `result.fun` returns the optimum function values.
- ② `result.x` returns the arguments corresponding to the optimum.
- ③ `result.success` returns True if optimization complete.

Numerical integration

Numerical integration is done via the `scipy.integrate` package.

```
1 import scipy.integrate as integr
```

There are several methods to numerically integrate a function (say $f(x) = \sin x + \frac{x}{2}$); fixed Gaussian quadrature, adaptive quadrature, Romberg integration....

All are approximations of the same thing, though...

```
1 integr.fixed_quad(f, lmin, lmax)[0]
2 integr.quad(f, lmin, lmax)[0]
3 integr.romberg(f, lmin, lmax)[0]
```
