

Basic Statistics pt1: Fitting

Carver J. Bierson

February 2, 2018

Caveat: In this class I will give you some simple statistical tools and teach you how to use them. This is not a statistics course so we will not go into much detail about when different measures are meaningful. Different statistical techniques and measure often have different assumptions built into them that we will only touch on. When applying these tools on real data make sure that you understand your statistical tools and they apply to your data.

Later in the course we will revisit many of these topics to discuss goodness of fit and choosing the correct model for your data.

1 Random numbers

When doing statistics and testing statistics functions it is helpful to be able to generate random data. For this class we will use the `numpy.random` module. This module contains functions to draw values from many different distributions. Below is a basic example of generating some random data.

```
1 import matplotlib.pyplot as plt
2
3 import numpy.random as rnd
4
5 set_uni=rnd.random(int(1E3)) # generate 1000 random uniformly
   distributed values
6
7 set_norm=rnd.randn(int(1E3)) # generate 1000 random normally
   distributed values
8
9
10 # Plot data
11 plt.figure()
12
13 Nbins=50 # number of bins in histogram
14
15 plt.subplot(121) # subplot for uniform data
16 plt.title('Uniform')
17 plt.hist(set_uni, Nbins, facecolor='green')
18
19 plt.subplot(122)
20 plt.title('Normal')
21 plt.hist(set_norm, Nbins, facecolor='red')
22
23 plt.show()
```

Listing 1: Create and plot a histogram of random points

1.1 Seeds

Computers can't generate truly random numbers. Instead they use sophisticated mathematical formulas that simulate random values. As an aside generating truly random values is a large area of research in computer science with applications in casinos and security. Most random number generation relies on a *seed* value. If this seed is the same, you will always get the same values. If you include

```
1 rnd.seed(1)
```

in the code above (before calling random) it will always produce the same 'random' numbers. Python seems to be good and using a changing quantity to supply this seed so you shouldn't need to worry about this too much (If you just run the same code twice you get different values). If you do find yourself in a case where you need a new seed using the least significant figure of the time is a common method.

2 Descriptive Statistics

Numpy gives us easy to use functions to calculate many of basic statistical properties of a data set. In addition to mean, median, and standard deviation numpy will also calculate data value quantiles and percentiles. A percentile gives the data value where X percent of the data is less than that value. For example the 25th percentile is the value that is larger than 25% of the data and the 50th percentile is the median. Percentiles are more robust to outliers than the standard deviation.

Below I use a random data set as an example.

```
1 import numpy as np
2 import numpy.random as rnd
3
4 set_norm=rnd.randn(int(1E3)) # generate 1000 random normally
   distributed values
5
6 print('The Mean is {:.2f}'.format(np.mean(set_norm))) # Print the
   data mean
7
8 print('The median is {:.2f}'.format(np.median(set_norm))) # Print
   the data median
9
10 print('The standard deviation is {:.2f}'.format(np.std(set_norm)))
    # Print the data standard deviation
11
12 print('The maximum value is {:.2f}'.format(np.max(set_norm))) #
   Print the maximum value
13
14
15 print('The 75th and 99th data percentile are {:.2f}, {:.2f}'.
   format(
16     np.percentile(set_norm,75),
17     np.percentile(set_norm,99))) # Print the data mean
```

Listing 2: Create and plot a histogram of random points

3 Fitting

Often with data we are interested in fitting some sort of model to that data and evaluating the quality of fit. Almost all model fitting is based on minimizing the squared error,

$$E = \sum_i (p(x_i) - y_i)^2 \quad (1)$$

where y_i is one of our data points and $p(x_i)$ is some model evaluated at the x value that corresponds to y_i . I will also show some basic error estimates for fits. Note that often the error estimates that are return are under-estimates of the real error.

4 Linear fits

Lets start with the simple example of trying to fit a line.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import numpy.random as rnd
4
5 # Generate some fake data
6 Datax=rnd.random(50)*10 # start with 50 uniform points
7
8 # Set the 'true' solution to y=3*x-2
9 Datay= 3*Datax-2
10
11 # add noise with a standard deviation of 2 to the data
12 Datay+=2*rnd.randn(50)
13
14
15 # find best linear fit
16 fit=np.polyfit(Datax,Datay,1) # 1 is the power of the fit , 3 would
    be a cubic fit
17
18 print(fit)
19
20 Plotx=np.linspace(0,10) # create an x array for plotting the fit
21 Ploty=np.polyval(fit,Plotx)
22
23 plt.figure()
24 plt.plot(Datax,Datay,'b.',label='Data') #plot data as points
25 plt.plot(Plotx,Ploty,'r-',label='Fit') #plot data as points
26
27 plt.legend()
```

In the example above I start with the equation $y = 3x - 2$. When I ran the code I got a fit of $y = 3.1x - 2.6$. We would like to put error estimates on those fit values to evaluate if we are really doing a good job. For this we simply need to return more values from polyfit. If we return the covariance matrix from polyfit, the diagonal elements of that matrix give the estimated standard deviation of each paramter. **cov matrix contains σ^2 , not σ . Should be corrected**

```
1 import numpy as np
2 import numpy.random as rnd
3
4 # Generate some fake data
5 Datax=rnd.random(50)*10 # start with 50 uniform points
6
```

```

7 # Set the 'true' solution to y=3*x-2
8 Datay= 3*Datax-2
9
10 # add noise with a standard deviation of 2 to the data
11 Datay+=2*rnd.randn(50)
12
13
14 # find best linear fit
15 fit , covariance=np.polyfit(Datax,Datay,1,cov=True) # return
    covariance matrix
16
17 print('Fitting model y=mx+b')
18 print('m={:0.2f}+/-{:0.2f}'.format(fit[0],2*covariance[0,0]))
19 print('b={:0.2f}+/-{:0.2f}'.format(fit[1],2*covariance[1,1]))

```

When I ran this code I got. Your results will vary a bit each time because of the random arrays

```

1 Fitting model y=mx+b
2 m=3.00+/-0.02
3 b=-2.39+/-0.61

```

4.1 Log space fitting

Many equations in science take the form

$$y = Ae^{bx} \quad (2)$$

or

$$y = Ax^b \quad (3)$$

Lets say we want to solve for the parameters A and b . We can solve both of these types of equations through a simple linearization. For both equation I will take the log of both sides and use some log algebra.

$$y = Ae^{bx} \quad (4)$$

$$\ln(y) = \ln(Ae^{bx}) \quad (5)$$

$$\ln(y) = \ln(A) + bx \quad (6)$$

$$(7)$$

and

$$y = Ax^b \quad (8)$$

$$\ln(y) = \ln(Ax^b) \quad (9)$$

$$\ln(y) = \ln(A) + b\ln(x) \quad (10)$$

Below is a code solving the second form.

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 import numpy.random as rnd
4
5 # Generate some fake data
6 Datax=1+rnd.random(50)*10 # start with 50 uniform points
7

```

```

8 # Set the 'true' solution to  $y=3*x^{2.5}$ 
9 Datay= 3*Datax**2.5
10
11 #move to logspace
12 lnDatay=np.log(Datay)
13 lnDatax=np.log(Datax)
14
15 # add noise in the logspace with a standard deviation of 0.5 to the
    data
16 lnDatay+=0.5*rnd.randn(50)
17
18
19 # find best linear fit
20 fit , covariance=np.polyfit(lnDatax ,lnDatay ,1 ,cov=True) # return
    covariance matrix
21
22 A=np.exp(fit [1])
23 b=fit [0]
24 Astd=np.exp(covariance [1,1])
25 bstd=covariance [0,0]
26
27 print('Fitting model  $y=Ax^b$ ')
28 print('A={:0.2f}+/-{:0.2f}'.format(A,2*Astd))
29 print('b={:0.2f}+/-{:0.2f}'.format(b,2*bstd))
30
31
32 Plotx=np.linspace(1,11) # create an x array for plotting the fit
33 Ploty=np.exp(np.polyval(fit ,np.log(Plotx)))
34
35 # Plot the data in log space
36 plt.figure()
37 plt.plot(lnDatax ,lnDatay , 'b.' , label='Data') #plot data as points in
    log space
38
39 plt.xlabel('ln(x)')
40 plt.ylabel('ln(y)')
41
42 plt.legend()
43
44 # plot in linear space
45 plt.figure()
46 plt.plot(Datax ,Datay , 'b.' , label='Data') #plot data as points
47 plt.plot(Plotx ,Ploty , 'r-' , label='Fit') #plot data as points
48
49 plt.xlabel('x')
50 plt.ylabel('y')
51
52 plt.legend()

```

When I ran this code I got.

```

1 Fitting model  $y=Ax^b$ 
2 A=3.69+/-2.09
3 b=2.41+/-0.03

```

Nicely this is consistent with my input model.

4.2 Non-linear fitting

Sometimes you want to fit a non-linear model to your data. Here there is no analytically solution to getting a fit. That said there are many numerical methods that can get these fits through iteration (think loops).

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 import numpy.random as rnd
4
5 def sinfunc(x, p0, p1, p2, p3):
6     return p0 + p1*np.sin(p2*(x-p3))
7
8 # Generate some fake data
9 Datax=rnd.random(50)*2*np.pi # start with 50 uniform points
10
11 # Set the 'true' solution to y=2*sin(x-0.1)
12 Datay= 2*np.sin(Datax-0.1)
13
14
15
16 # add noise with a standard deviation of 0.25 to the data
17 Datay+=0.25*rnd.randn(50)
18
19 from scipy.optimize import curve_fit # load non-linear fitting
    function
20 p, pcov = curve_fit(sinfunc, Datax, Datay) # perform fit
21
22 # p contains the fit
23 # pcov contains the covariance matrix
24
25 # print result
26 print('Fitting model y=p0 + p1*sin(p2*(x-p3))')
27 print('p0={:0.2f}+/-{:0.2f}'.format(p[0],2*pcov[0,0]))
28 print('p1={:0.2f}+/-{:0.2f}'.format(p[1],2*pcov[1,1]))
29 print('p2={:0.2f}+/-{:0.2f}'.format(p[2],2*pcov[2,2]))
30 print('p3={:0.2f}+/-{:0.2f}'.format(p[3],2*pcov[3,3]))
31
32 Plotx=np.linspace(0,2*np.pi) # create an x array for plotting the
    fit
33 Ploty=sinfunc(Plotx,p[0],p[1],p[2],p[3])
34
35 # plot Solution
36 plt.figure()
37 plt.plot(Datax,Datay,'b.',label='Data') #plot data as points
38 plt.plot(Plotx,Ploty,'r-',label='Fit') #plot data as points
39
40 plt.xlabel('x')
41 plt.ylabel('y')
42
43 plt.legend()

```

When I ran this code I got.

```

1 Fitting model y=p0 + p1*sin(p2*(x-p3))
2 p0=0.02+/-0.00
3 p1=2.05+/-0.01
4 p2=1.00+/-0.00
5 p3=0.05+/-0.01

```

The fit values are fairly good but the error estimates are clearly too low. This is not uncommon for fitting routines so be careful how you interpret your results (and other peoples).

4.3 Spline fitting and interpolation

Sometimes you are less interested in fitting a physical model and just need to interpolate between data points. For this spline curves are often a good solution.

In the example below I sparsely sample a sin curve, then compare a linear and spline interpolation.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4
5 # Sample a sin curve
6 Datax=np.linspace(0,2*np.pi,4) # start with 50 uniform points
7
8 # Create a complex y data set
9 Datay= 2*np.sin(Datax-0.1)+2*Datax
10
11
12 import scipy.interpolate as interp
13
14 Plotx=np.linspace(0,2*np.pi,100) # create an x array for plotting
    the fit
15
16 Lineary=np.interp(Plotx,Datax,Datay)
17 Spliney=interp.spline(Datax,Datay,Plotx)
18 Truey= 2*np.sin(Plotx-0.1)+2*Plotx
19
20
21 # plot Solution
22 plt.figure()
23
24 plt.plot(Plotx,Truey,'b-', label='True Fit')
25 plt.plot(Plotx,Lineary,'r--', label='Linear Fit')
26 plt.plot(Plotx,Spliney,'g--', label='Spline Fit')
27
28 plt.plot(Datax,Datay,'b.',label='Data') #plot data as points
29
30
31 plt.xlabel('x')
32 plt.ylabel('y')
33
34 plt.legend()
```

Note that spline fits tend to assume things are very smooth. If your data is not smooth this is not going to give you a good interpolation!