

Least squares fit to a straight line

Derivation of the least-squares formulae for slope and intercept

Suppose we want to fit a dependent data set y_i to an independent set x_i according to the model

$$y = mx + b$$

In the following, we will assume that all of the y values *have the same uncertainty* measured by a standard deviation σ_y .

For trial values of slope m and intercept b we'd like to minimize the difference between the predicted value y from the measured value y_i for all of the values of x_i . We call these differences the "residuals" and define a measure χ^2 for the goodness of fit for a given choice of parameters:

$$\chi^2 = \sum_{i=1}^N \frac{1}{\sigma_{yi}^2} [y_i - (mx_i + b)]^2$$

Note that this is a weighted sum of the squares of the residuals divided by the squares of data uncertainties σ_{yi} that depend on which data value y_i is used.

We get the best fit for a minimum of the above quantity (hence the name "least squares"), which requires the derivatives with respect to m and b to vanish.

$$\frac{\partial}{\partial m} \chi^2 = 0.$$

$$\frac{\partial}{\partial b} \chi^2 = 0.$$

Applying these derivatives to the definition of χ^2 :

$$\sum_{i=1}^N \frac{1}{\sigma_{yi}^2} 2[y_i - (mx_i + b)](-x_i) = 0.$$

$$\sum_{i=1}^N \frac{1}{\sigma_{yi}^2} 2[y_i - (mx_i + b)](-1) = 0.$$

Rearranging, we have two equations in two unknowns m and b :

$$\left(\sum_{i=1}^N \frac{x_i^2}{\sigma_{yi}^2} \right) m + \left(\sum_{i=1}^N \frac{x_i}{\sigma_{yi}^2} \right) b = \sum_{i=1}^N \frac{x_i y_i}{\sigma_{yi}^2}.$$

$$\left(\sum_{i=1}^N \frac{x_i}{\sigma_{yi}^2} \right) m + \left(\sum_{i=1}^N \frac{1}{\sigma_{yi}^2} \right) b = \sum_{i=1}^N \frac{y_i}{\sigma_{yi}^2}.$$

Define weighted averages as follows:

$$\langle x \rangle \equiv \frac{\tilde{\sigma}_y^2}{N} \sum_{i=1}^N \frac{x_i}{\sigma_{yi}^2} \left(= \frac{1}{N} \sum_{i=1}^N x_i \text{ if all } \sigma_{yi} \text{ are equal} \right).$$

$$\langle y \rangle \equiv \frac{\tilde{\sigma}_y^2}{N} \sum_{i=1}^N \frac{y_i}{\sigma_{yi}^2} \left(= \frac{1}{N} \sum_{i=1}^N y_i \text{ if all } \sigma_{yi} \text{ are equal} \right).$$

$$\langle x^2 \rangle \equiv \frac{\tilde{\sigma}_y^2}{N} \sum_{i=1}^N \frac{x_i^2}{\sigma_{yi}^2} \left(= \frac{1}{N} \sum_{i=1}^N x_i^2 \text{ if all } \sigma_{yi} \text{ are equal} \right).$$

$$\langle xy \rangle \equiv \frac{\tilde{\sigma}_y^2}{N} \sum_{i=1}^N \frac{x_i y_i}{\sigma_{yi}^2} \left(= \frac{1}{N} \sum_{i=1}^N x_i y_i \text{ if all } \sigma_{yi} \text{ are equal} \right).$$

Here

$$\frac{1}{\tilde{\sigma}_y^2} \equiv \frac{1}{N} \sum_{i=1}^N \frac{1}{\sigma_{yi}^2} \left(= \frac{1}{\sigma_y^2} \text{ if all } \sigma_{yi} \text{ are equal to the same value } \sigma_y \right).$$

$$\tilde{\sigma}_y^2 \equiv \frac{N}{\sum_{i=1}^N \frac{1}{\sigma_{yi}^2}} \left(= \sigma_y^2 \text{ if all } \sigma_{yi} \text{ are equal to the same value } \sigma_y \right).$$

After multiplying through by a common factor of $\tilde{\sigma}_y^2$, the equations for slope and intercept become

$$\langle x^2 \rangle m + \langle x \rangle b = \langle xy \rangle.$$

$$\langle x \rangle m + b = \langle y \rangle.$$

Using this notation, solve the second of the simultaneous equations for b :

$$b = \langle y \rangle - m \langle x \rangle.$$

This seems to be an intuitively plausible result.

Substitute this value for b into the first equation to get:

$$\langle x^2 \rangle m + \langle x \rangle [\langle y \rangle - m \langle x \rangle] = \langle xy \rangle.$$

Solve for m :

$$m = \frac{\langle xy \rangle - \langle x \rangle \langle y \rangle}{\langle x^2 \rangle - \langle x \rangle^2}.$$

Once we have m , we simply use this in the expression above to calculate b :

$$b = \langle y \rangle - m \langle x \rangle.$$

Uncertainties in the slope and intercept estimates

*The following results can be derived through propagation of errors in the expressions for m and b above - after explicitly writing out the averages as sums and taking derivatives with respect to the uncertain quantities y_i .

(1) The estimated uncertainty in the slope is given by

$$\sigma_{m_{est}} = \frac{\tilde{\sigma}_y}{\sqrt{N}} \sqrt{\frac{1}{\langle x^2 \rangle - \langle x \rangle^2}}.$$

(2) The estimated uncertainty in the intercept is given by

$$\sigma_{b_{est}} = \frac{\tilde{\sigma}_y}{\sqrt{N}} \sqrt{\frac{\langle x^2 \rangle}{\langle x^2 \rangle - \langle x \rangle^2}} = \sigma_{m_{est}} \sqrt{\langle x^2 \rangle}.$$

An interesting feature of these results is that if, from one experimental run to the next, the same set of values of the independent variable x are used, then the uncertainties in slope and intercept will only vary because of variations in the set of *uncertainties* of y , not because of variations in the y values themselves. This is demonstrated in the numerical exploration below when doing weighted fits with specified uncertainties in y .

Again, if all of the uncertainties in y are assumed to be the same value σ_y , then we replace $\tilde{\sigma}_y$ with σ_y in the above expressions. If the uncertainty σ_y is not known *a priori* then it can be estimated by

$$\sigma_{y_{est}} \approx \sqrt{\frac{1}{N-2} \sum_{i=1}^N [y_i - (mx_i + b)]^2},$$

where m and b are the estimated values of slope and intercept. When this estimate is used, the uncertainties in slope and intercept will vary from run to run as the set of values of y used in the estimate varies.

References

<https://www.asc.ohio-state.edu/gan.1/teaching/spring04/Chapter7.pdf> (<https://www.asc.ohio-state.edu/gan.1/teaching/spring04/Chapter7.pdf>)

Reference: see P. R. Bevington & D. K. Robinson, *Data Reduction and Error Analysis for the Physical Sciences*, 3rd ed., (McGraw-Hill, 2003) Chapter 6.

Code to do linear regression in several ways

The following approaches are used to do a linear regression for a small data set.

If you want to use the code, edit the variables:

```
my-name
xdata (an array)
ydata (an array)
sigma_ydata (an array of uncertainty estimates for variable y)
```

The code then calculates fitting parameters m and b for the linear model $y = mx + b$ along with uncertainties in the slope using the following methods:

- explicit calculation using formulas from the above derivations but without weighting by the individual point variance estimates
- same as above but with weighting
- unweighted fit using the module `linregress` from the `scipy.stats` library
- unweighted and weighted fits using the module `curve_fit` from the `scipy.optimize` module

You should see that the unweighted results all agree and separately how the weighted results agree.

The best fit lines (weighted and unweighted) are then plotted through the data points, which have error-bars.

In [15]:

```
1  '''
2  multLS.py
3  Performs least squares fit multiple ways, both unweighted and weighted
4  Randall Tagg 07-Oct-2020
5  '''
6
7  import math
8  import numpy as np
9  import matplotlib.pyplot as plt
10 from scipy.stats import linregress
11 from scipy.optimize import curve_fit
12 import datetime
13
14 ##### EDIT NAME AND DATA ARRAYS HERE #####
```

```

15
16 my_name = "Randall Tagg"
17
18 #Enter x data
19 xdata = np.asarray([0.,1.,2.,3.,4.,5.,6.])
20 Ndata = xdata.size
21
22 #Enter y data
23 ydata = np.asarray([4.111,7.632,9.563,14.530,16.269,19.003,20.752])
24
25 #Enter standard deviations of y data
26 sigma_ydata = np.asarray([0.5,0.7,0.3,1.0,0.4,0.2,0.6])
27 #Uncomment the following two lines and edit if needed for
28 # uniform weighting
29 #sigmaScaleFactor = 1.
30 #sigma_ydata = sigmaScaleFactor*np.asarray([1.,1.,1.,1.,1.,1.,1.])
31
32
33 # Create a time stamp for each run
34 time_stamp_format = "%Y-%m-%d %H:%M:%S"
35 print(my_name,datetime.datetime.now().strftime(time_stamp_format),'\
36 print('x ',xdata)
37 print('y',ydata)
38 print('sigma_ydata',sigma_ydata)
39
40 #####
41
42
43 def flinear(xx,mm,bb): # define the linear function (for use with t
44     return mm*xx + bb
45
46
47 #####
48
49
50 # Explicitly compute the fit parameters and their uncertainties
51 # with no weighting
52
53 print('\nCompute fit parameters explicitly from derived equations')
54 xav = 0
55 yav = 0
56 xyav = 0
57 x2av = 0
58 for i in range(0,Ndata):
59     xav = xav+xdata[i]/Ndata
60     yav = yav+ydata[i]/Ndata
61     x2av = x2av+xdata[i]*xdata[i]/Ndata
62     xyav = xyav+xdata[i]*ydata[i]/Ndata
63
64 m_fit = (xyav-xav*yav)/(x2av-xav**2)
65 b_fit = yav - m_fit*xav
66
67 ypredict = np.zeros(Ndata)
68 variance_fit = 0

```

```

69 for i in range(0,Ndata):
70     ypredict[i] = flinear(xdata[i],m_fit,b_fit)
71     # In the following, 2 is subtracted from Ndata because 2 degrees
72     # freedom (slope and intercept) were used to estimate sigma_y
73     variance_fit = variance_fit + (ydata[i]-ypredict[i])**2
74 sigma_y_estimate = math.sqrt(variance_fit/(Ndata-2))
75 sigma_m_fit = sigma_y_estimate/math.sqrt(Ndata*(x2av-xav**2))
76 sigma_b_fit = sigma_m_fit*math.sqrt(x2av)
77
78 print('Unweighted   m_fit = {0:0.3f} \u00B1 {1:0.3f}   b_fit = {2:0
79       .format(m_fit,sigma_m_fit,b_fit,sigma_b_fit))
80
81
82 #####
83
84
85 # Now explicitly compute the fit parameters and their uncertainties
86 # using variance weighting
87
88 xav = 0
89 yav = 0
90 xyav = 0
91 x2av = 0
92 sum_one_over_sigma_y_squared = 0
93 for i in range(0,Ndata):
94     xav = xav+xdata[i]/(sigma_ydata[i])**2
95     yav = yav+ydata[i]/(sigma_ydata[i])**2
96     x2av = x2av+xdata[i]*xdata[i]/(sigma_ydata[i])**2
97     xyav = xyav+xdata[i]*ydata[i]/(sigma_ydata[i])**2
98     sum_one_over_sigma_y_squared = sum_one_over_sigma_y_squared \
99     + 1/(sigma_ydata[i])**2
100
101 xav = xav/sum_one_over_sigma_y_squared
102 yav = yav/sum_one_over_sigma_y_squared
103 x2av = x2av/sum_one_over_sigma_y_squared
104 xyav = xyav/sum_one_over_sigma_y_squared
105 sigma_y_tilde = math.sqrt(Ndata/sum_one_over_sigma_y_squared)
106
107 m_fit_w = (xyav-xav*yav)/(x2av-xav**2)
108 b_fit_w = yav - m_fit_w*xav
109
110 sigma_m_fit_w = sigma_y_tilde/math.sqrt(Ndata*(x2av-xav**2))
111 sigma_b_fit_w = sigma_m_fit_w*math.sqrt(x2av)
112
113 print('Weighted   m_fit_w = {0:0.3f} \u00B1 {1:0.3f}   b_fit_w = {2:0
114       .format(m_fit_w,sigma_m_fit_w,b_fit_w,sigma_b_fit_w ))
115
116
117 #####
118
119
120 # Now use scipy.stats linregress module instead.
121 # https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.l
122 print()

```

```

123 m_stats, b_stats, r_value, p_value, std_err = linregress(xdata,ydata)
124 print('Now use numpy.stats linregress module p_value = {0:0.2e} r_value = {0:0.2e}'.format(p_value, r_value))
125 print('m_stats = {0:0.3f} \u00B1 {1:0.3f}      b_stats={2:0.3f}'.format(m_stats, r_value, b_stats))
126
127
128 #####
129
130
131 # Finally, use scipy.optimize curve_fit module, first unweighted and
132 # https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.curve_fit.html
133 # See also
134 # https://scipython.com/book/chapter-8-scipy/examples/weighted-and-unweighted-linear-fitting/
135 print('\nFinally use numpy.optimize curve_fit module')
136 popt, pcov = curve_fit(flinear, xdata, ydata)
137 m_opt,b_opt = popt
138 sigma_m_opt,sigma_b_opt = np.sqrt(np.diag(pcov))
139 print('m_opt = {0:0.3f} \u00B1 {1:0.3f}      b_opt = {2:0.3f} \u00B1 {3:0.3f}'.format(m_opt, sigma_m_opt, b_opt, sigma_b_opt))
140
141 popt_w, pcov_w = curve_fit(flinear, xdata, ydata, p0=None, sigma=sigma_y_tilde)
142 m_opt_w,b_opt_w = popt_w
143 sigma_m_opt_w,sigma_b_opt_w = np.sqrt(np.diag(pcov_w))
144 print('m_opt_w = {0:0.3f} \u00B1 {1:0.3f}      b_opt_w = {2:0.3f} \u00B1 {3:0.3f}'.format(m_opt_w, sigma_m_opt_w, b_opt_w, sigma_b_opt_w))
145
146
147
148 #####
149
150
151 print()
152 print('sigma_y using fit = {0:0.3f}'.format(sigma_y_estimate))
153 print('sigma_y_tilde = {0:0.3f}'.format(sigma_y_tilde))
154
155
156 #####
157
158
159 # This is needed for plotting the line fitted using variance weighted curve_fit
160 ypredict_w = np.zeros(Ndata)
161 for i in range(0,Ndata):
162     ypredict_w[i] = flinear(xdata[i],m_opt_w,b_opt_w)
163     # ypredict_w[i] = m_opt_w*xdata[i]+b_opt_w
164
165
166 #####
167
168
169 # Now plot the data and best fit lines for unweighted and weighted curve_fit
170 datapoints, = plt.plot(xdata,ydata,'r.',markersize=12,label='data')
171 plt.errorbar(xdata,ydata,yerr=sigma_ydata,marker='.',linestyle = 'none')
172 fittedline, = plt.plot(xdata,ypredict,label='fitted line')
173 fittedline_w, = plt.plot(xdata,ypredict_w,label='fitted line, weighted')
174 plt.xlabel('x')
175 plt.ylabel('y')
176 plt.title('Linear fit to data')

```

```

176 plt.title('Linear fit to data',
177 plt.legend(handles=[datapoints,fittedline,fittedlinew])
178 plt.show()
179

```

Randall Tagg 2021-12-04 16:30:38

```

x [0. 1. 2. 3. 4. 5. 6.]
y [ 4.111  7.632  9.563 14.53  16.269 19.003 20.752]
sigma_ydata [0.5 0.7 0.3 1.  0.4 0.2 0.6]

```

Compute fit parameters explicitly from derived equations

Unweighted $m_{\text{fit}} = 2.835 \pm 0.163$ $b_{\text{fit}} = 4.619 \pm 0.586$

Weighted $m_{\text{fit}_w} = 2.967 \pm 0.079$ $b_{\text{fit}_w} = 4.059 \pm 0.325$

Now use numpy.stats linregress module $p_value = 1.14e-05$ $r_value=0.992$

$m_stats = 2.835 \pm 0.163$ $b_stats=4.619$

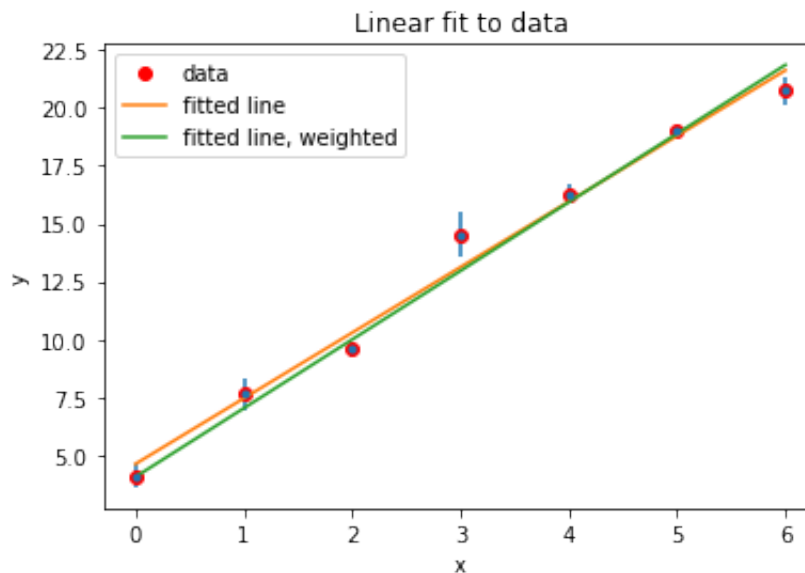
Finally use numpy.optimize curve_fit module

$m_{\text{opt}} = 2.835 \pm 0.163$ $b_{\text{opt}} = 4.619 \pm 0.586$

$m_{\text{opt}_w} = 2.967 \pm 0.079$ $b_{\text{opt}_w} = 4.059 \pm 0.325$

σ_y using fit = 0.861

$\sigma_{y_tilde} = 0.366$



Code to read in data from a file and fit to a straight line

The following will do a weighted fit using the `curve_fit` module from the `scipy.optimize` library for data read from a comma delimited (.csv) file. There should be two or three columns:

- xdata
- ydata
- uncertainty estimates for ydata If the third column is missing then uniform weighting will be used and an estimate of the y variance will be made using the residuals (i.e. the differences between best fit and data).

The first line of the file should contain the variable names. The second line of the file should contain the units.

Here is an example:

```
time,position,sigma
s,m,m
0,15,1.2
1,17.1,0.5
2,18.9,1.1
3,21.3,1.2
4,22.7,1.1
```

This renders in table form as

time	position	sigma
s	m	m
0	15	1.2
1	17.1	0.5
2	18.9	1.1
3	21.3	1.2
4	22.7	1.1

Edit the variables in the code below before running:

```
myname = (your name for the output time stamp)
myDataFile = (name of the data file, which should have extension
.csv)
plotTitle = (a text string)
```

It is likely also to be necessary to edit the data formatting in the print statements. (The example uses the format 0.3f for all values.)

```

In [83]: 1 %matplotlib inline
2 '''Code to read in data including uncertainties and do linear least s
3
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import math,datetime
7 from scipy.optimize import curve_fit
8
9 def flinear(xx,mm,bb): # define the linear function (for use with th
10     return mm*xx + bb
11
12 # *****USER EDITED PORTION *****
13
14 #Put in the name of the desired data file in the following line
15 #and edit the x and y labels to suit the paritcular data
16 my_name = "J.J. Datasmith"
17 myDataFile = 'testData.csv'
18 plotTitle = 'Least Squares Fit'
19
20 # *****
21
22 #Read in the header
23 f = open(myDataFile)
24 varnames = f.readline()
25 units = f.readline()
26
27 #Read in the data
28 mydata = np.loadtxt(myDataFile,delimiter=',',skiprows=2)
29
30 if (mydata.shape[1] == 3):
31     difvar = True
32     xname,yname,signame = varnames.split(',')
33     xunit,yunit,sigunit = units.split(',')
34 else:
35     difvar = False
36     xname,yname = varnames.split(',')
37     yname = yname.rstrip(yname[-1]) # remove the line feed
38     xunit,yunit = units.split(',')
39     yunit = yunit.rstrip(yunit[-1]) # remove the line feed
40
41 #print(difvar)
42 #print(mydata.shape[1])
43
44 x = mydata[:,0] # x data
45 y = mydata[:,1] # y data
46 if (difvar):
47     sigma_y = mydata[:,2] # uncertainties in y data
48 else:
49     sigma_y = np.ones(Ndata) # equal weighting
50
51 # Use numpy.optimize curve_fit module to make a least squares fit of
52 # data to a straight line, using variance weighting (hence the _w des
53 popt_w, pcov_w = curve_fit(flinear, x, y, p0=None, sigma=sigma_y, abs
54 m,b = popt_w

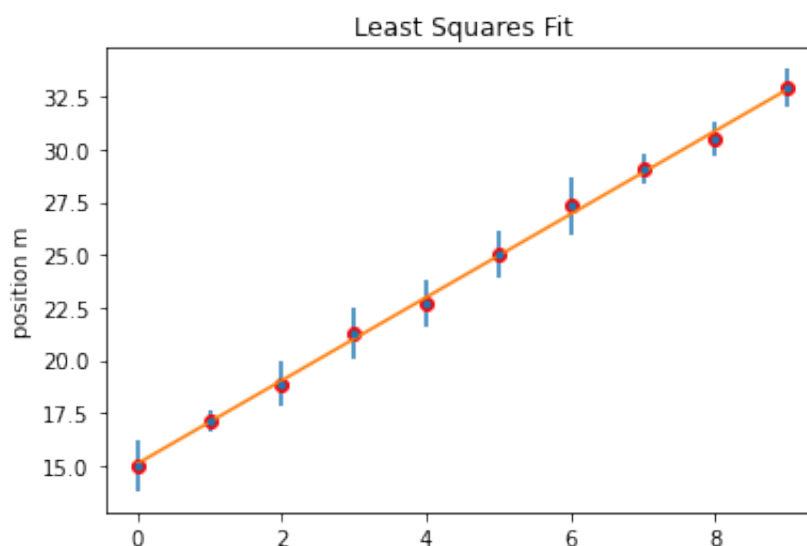
```

```

55 sigma_m, sigma_b = np.sqrt(np.diag(pcov_w))
56
57 #Compute y values predicted from best fit line
58 Ndata = y.size
59 ypredicted = np.zeros(Ndata)
60 variance_fit = 0
61 tempsig = 0
62 for i in range(0,Ndata):
63     ypredicted[i] = flinear(x[i],m,b)
64     variance_fit = variance_fit + (y[i]-ypredicted[i])**2
65     if(difvar): tempsig = tempsig + 1/sigma_y[i]**2
66 # In the following, 2 is subtracted from Ndata because 2 degrees of
67 # freedom (slope and intercept) were used to estimate sigma_y
68 sigma_y_estimate = math.sqrt(variance_fit/(Ndata-2))
69 if(difvar): sigma_y_tilde = math.sqrt(Ndata/tempsig)
70
71 # Create a time stamp for each run
72 time_stamp_format = "%Y-%m-%d %H:%M:%S"
73 print(my_name,datetime.datetime.now().strftime(time_stamp_format),'\n')
74
75 plt.plot(x,y,'r.',markersize=12,label='data')
76 if(difvar): plt.errorbar(x,y,yerr=sigma_y,marker='.',linestyle = 'none')
77 plt.plot(x,ypredicted,label='best fit line')
78 xlabel = xname + ' ' + xunit
79 plt.xlabel(xlabel)
80 ylabel = yname + ' ' + yunit
81 plt.ylabel(ylabel)
82 plt.title(plotTitle)
83
84 plt.show()
85
86 print('Slope m = {0:0.3f} \u00B1 {1:0.3f}'.format(m,sigma_m),yunit,'/')
87 print('Intercept b = {0:0.3f} \u00B1 {1:0.3f}'.format(b,sigma_b),yunit)
88 print('Fit estimated y std dev = {0:0.3f}'.format(sigma_y_estimate),yunit)
89 if(difvar): print('Inverse avg y std dev = {0:0.3f}'.format(sigma_y_tilde),yunit)

```

J.J. Datasmith 2023-03-03 19:22:04



time s

Slope $m = 1.969 \pm 0.089 \text{ m / s}$
Intercept $b = 15.113 \pm 0.466 \text{ m}$
Fit estimated y std dev = 0.255 m
Inverse avg y std dev = 0.870 m

Code to explore the effects size of randomness and other details

The following creates a synthetic data set with adjustable uncertainties in order to explore how these affect the fit.

In [32]:

```
1 '''
2 exploreLSv2.py
3 Code to explore least squares fitting: an artificial data set
4 is generated using a random number generator that uses a Gaussian
5 distribution of errors with standard deviation specified by the
6 parameter sigma. The effect of changing the size of sigma on the
7 quality of fit can be explored. Both unweighted and weighted fitting
8 is done.
9 Randall Tagg 07-Oct-2020
10 '''
11
12 import math
13 import numpy as np
14 import matplotlib.pyplot as plt
15 from scipy.stats import linregress
16 from scipy.optimize import curve_fit
17 import datetime
18
19 # Create a time stamp for each run
20 my_name = "Randall Tagg"
21 time_stamp_format = "%Y-%m-%d %H:%M:%S"
22 print(my_name,datetime.datetime.now().strftime(time_stamp_format),'\n')
23
24 def flinear(xx,mm,bb): # define the linear function (for use with the
25     return mm*xx + bb
26
27 #Actual slope and intercept for exact line
28 m=3.
29 b=4.
30
31
32 #Create synthetic data for exploring how fits depend on the distribution
33
34 xdata = np.asarray([0.,1.,2.,3.,4.,5.,6.])
35 Ndata = xdata.size
36
37 # Uncertainties to use in generating data using a normal distribution
38 sigma_scale_factor = 2
```

```

38 sigma_scale_factor = 2.
39 sigma_ydata = sigma_scale_factor*np.asarray([0.5,0.7,0.3,1.0,0.4,0.2,
40
41 #sigma_ydata = 2*np.asarray([0.5,0.5,0.5,0.5,0.5,0.5,0.5]) # use this
42
43 yexact = np.zeros(Ndata)
44 ydata = np.zeros(Ndata)
45 ydiff = np.zeros(Ndata)
46 ydiff_sum = 0
47 #Synthesize the y data
48 for i in range(0,Ndata):
49     yexact[i] = flinear(xdata[i],m,b)
50     ydata[i] = np.random.normal(yexact[i], sigma_ydata[i], 1) # create
51     ydiff[i] = ydata[i]-yexact[i]
52     ydiff_sum = ydiff_sum + ydiff[i]
53 ydiff_mean = ydiff_sum/Ndata
54 #Comment out the following line if using empirical mean for ydiff; ot
55 ydiff_mean = 0 # override previous result to force assumed zero mean
56
57 yvariance_sum = 0
58 for i in range(0,Ndata):
59     yvariance_sum = yvariance_sum + (ydiff[i]-ydiff_mean)**2
60 sigma_y_sample = math.sqrt(yvariance_sum/(Ndata-1))
61
62 print('x ',xdata)
63 print('Exact y ',yexact)
64 print('Random y [{0:.3f} {1:.3f} {2:.3f} {3:.3f} {4:.3f} {5:.3f} {6:.3f}]'.
65       .format(ydata[0],ydata[1],ydata[2],ydata[3],ydata[4],
66               ydata[5],ydata[6]))
67 print('sigma_ydata [{0:.3f} {1:.3f} {2:.3f} {3:.3f} {4:.3f} {5:.3f} {6:.3f}]'.
68       .format(sigma_ydata[0],sigma_ydata[1],sigma_ydata[2],sigma_ydata[3],
69               sigma_ydata[4],sigma_ydata[5],sigma_ydata[6]))
70
71 print('\nm_exact = {0:.3f}    b_exact = {1:.3f}'.format(m, b))
72
73
74 # Explicitly compute the fit parameters and their uncertainties with
75
76 print('\nCompute fit parameters explicitly from derived equations')
77 xav = 0
78 yav = 0
79 xyav = 0
80 x2av = 0
81 for i in range(0,Ndata):
82     xav = xav+xdata[i]/Ndata
83     yav = yav+ydata[i]/Ndata
84     x2av = x2av+xdata[i]*xdata[i]/Ndata
85     xyav = xyav+xdata[i]*ydata[i]/Ndata
86
87 m_fit = (xyav-xav*yav)/(x2av-xav**2)
88 b_fit = yav - m_fit*xav
89
90 ypredict = np.zeros(Ndata)
91 variance_fit = 0
92 for i in range(0,Ndata):

```

```

92 for i in range(0,Ndata):
93     ypredict[i] = flinear(xdata[i],m_fit,b_fit)
94     # In the following, 2 is subtracted from Ndata because 2 degrees
95     # freedom (slope and intercept) were used to estimate sigma_y
96     variance_fit = variance_fit + (ydata[i]-ypredict[i])**2
97 sigma_y_estimate = math.sqrt(variance_fit/(Ndata-2))
98 sigma_m_fit = sigma_y_estimate/math.sqrt(Ndata*(x2av-xav**2))
99 sigma_b_fit = sigma_m_fit*math.sqrt(x2av)
100
101 print('m_fit = {0:0.3f} \u00B1 {1:0.3f}      b_fit = {2:0.3f} \u00B1 {3
102       .format(m_fit,sigma_m_fit,b_fit,sigma_b_fit))
103
104
105 # Now explicitly computer the fit parameters and their uncertainties
106
107 xav = 0
108 yav = 0
109 xyav = 0
110 x2av = 0
111 sum_one_over_sigma_y_squared = 0
112 for i in range(0,Ndata):
113     xav = xav+xdata[i]/(sigma_ydata[i])**2
114     yav = yav+ydata[i]/(sigma_ydata[i])**2
115     x2av = x2av+xdata[i]*xdata[i]/(sigma_ydata[i])**2
116     xyav = xyav+xdata[i]*ydata[i]/(sigma_ydata[i])**2
117     sum_one_over_sigma_y_squared = sum_one_over_sigma_y_squared + 1/(
118
119 xav = xav/sum_one_over_sigma_y_squared
120 yav = yav/sum_one_over_sigma_y_squared
121 x2av = x2av/sum_one_over_sigma_y_squared
122 xyav = xyav/sum_one_over_sigma_y_squared
123 sigma_y_tilde = math.sqrt(Ndata/sum_one_over_sigma_y_squared)
124
125 m_fit_w = (xyav-xav*yav)/(x2av-xav**2)
126 b_fit_w = yav - m_fit_w*xav
127
128 sigma_m_fit_w = sigma_y_tilde/math.sqrt(Ndata*(x2av-xav**2))
129 sigma_b_fit_w = sigma_m_fit_w*math.sqrt(x2av)
130
131 print('m_fit_w = {0:0.3f} \u00B1 {1:0.3f}      b_fit_w = {2:0.3f} \u00B1 {3
132       .format(m_fit_w,sigma_m_fit_w,b_fit_w,sigma_b_fit_w ))
133
134
135
136 # Now use scipy.stats linregress module instead.
137 # https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.lin
138 print()
139 m_stats, b_stats, r_value, p_value, std_err = linregress(xdata,ydata)
140 print('Now use numpy.stats linregress module p_value = {0:0.2e}  r_va
141 print('m_stats={0:0.3f} \u00B1 {1:0.3f}      b_stats={2:0.3f}'.format(m
142
143
144
145 # Finally, use scipy.optimize curve_fit module, first unweighted and

```

```

146 # https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize
147 # See also
148 # https://scipython.com/book/chapter-8-scipy/examples/weighted-and-no
149 print('\nFinally use numpy.optimize curve_fit module')
150 popt, pcov = curve_fit(flinear, xdata, ydata)
151 m_opt, b_opt = popt
152 sigma_m_opt, sigma_b_opt = np.sqrt(np.diag(pcov))
153 print('m_opt = {0:0.3f} \u00B1 {1:0.3f}      b_opt = {2:0.3f} \u00B1 {3:0.3f}'
154       .format(m_opt, sigma_m_opt, b_opt, sigma_b_opt))
155 popt_w, pcov_w = curve_fit(flinear, xdata, ydata, p0=None, sigma=sigma_ydata)
156 m_opt_w, b_opt_w = popt_w
157 sigma_m_opt_w, sigma_b_opt_w = np.sqrt(np.diag(pcov_w))
158 print('m_opt_w = {0:0.3f} \u00B1 {1:0.3f}      b_opt_w = {2:0.3f} \u00B1 {3:0.3f}'
159       .format(m_opt_w, sigma_m_opt_w, b_opt_w, sigma_b_opt_w))
160
161
162
163 print()
164 print('sigma_y using diff from exact y values (zero mean diff) = {0:0.3f}'
165       .format(sigma_y_estimate))
166 #print('sigma_y using diff from exact y values (empirical mean diff) = {0:0.3f}'
167       .format(sigma_y_tilde))
168
169
170
171 # This is needed for plotting the line fitted using variance weighting
172 ypredict_w = np.zeros(Ndata)
173 for i in range(0, Ndata):
174     ypredict_w[i] = flinear(xdata[i], m_opt_w, b_opt_w)
175     # ypredict_w[i] = m_opt_w*xdata[i]+b_opt_w
176
177 datapoints, = plt.plot(xdata, ydata, 'r.', markersize=12, label='data')
178 plt.errorbar(xdata, ydata, yerr=sigma_ydata, marker='.', linestyle='none', label='exact')
179 exactline, = plt.plot(xdata, yexact, '--', label='exact line')
180 fittedline, = plt.plot(xdata, ypredict, label='fitted line')
181 fittedlinew, = plt.plot(xdata, ypredict_w, label='fitted line, weighted')
182 plt.xlabel('x')
183 plt.ylabel('y')
184 plt.title('Linear fit to data')
185 plt.legend(handles=[datapoints, exactline, fittedline, fittedlinew])
186 plt.show()
187

```

Randall Tagg 2021-12-04 17:32:21

```

x [0. 1. 2. 3. 4. 5. 6.]
Exact y [ 4.  7. 10. 13. 16. 19. 22.]
Random y [4.252 9.265 10.193 16.549 16.003 19.867 21.475]
sigma_ydata [1.000 1.400 0.600 2.000 0.800 0.400 1.200]

```

```
m_exact = 3.000    b_exact = 4.000
```

```

Compute fit parameters explicitly from derived equations
m_fit = 2.810 ± 0.288    b_fit = 5.513 ± 1.038

```

$m_{\text{fit}_w} = 3.012 \pm 0.158$ $b_{\text{fit}_w} = 4.559 \pm 0.650$

Now use numpy.stats.linregress module $p_value = 1.92e-04$ $r_value=0.975$

$m_{\text{stats}}=2.810 \pm 0.288$ $b_{\text{stats}}=5.513$

Finally use numpy.optimize.curve_fit module

$m_{\text{opt}} = 2.810 \pm 0.288$ $b_{\text{opt}} = 5.513 \pm 1.038$

$m_{\text{opt}_w} = 3.012 \pm 0.158$ $b_{\text{opt}_w} = 4.559 \pm 0.650$

σ_y using diff from exact y values (zero mean diff) = 1.773

σ_y using fit = 1.524

$\sigma_{y_tilde} = 0.733$

