

# Introduction to the Least Squares Fit

## Table of Contents

- 1. Uses for a Least Squares Fit: Linear Dependence**
- 2. Methods of Finding the Best Fit Line: Estimating, Using Excel, and Calculating Analytically**
- 3. Calculating a Least Squares Fit**
- 4. Uncertainty in the Dependent Variable, Slope, and Intercept**
- 5. Calculating the  $r^2$  Value**
- 6. Sample Data and Setting up the Spreadsheet**
  - 6.1 Example: Entering a Sum**
  - 6.2 Example: Entering a Formula**
  - 6.3 Formulae for Excel Using Sample Data**
  - 6.4 Using Excel's LINEST Function**

## 1. Uses for a Least Squares Fit: Linear Dependence

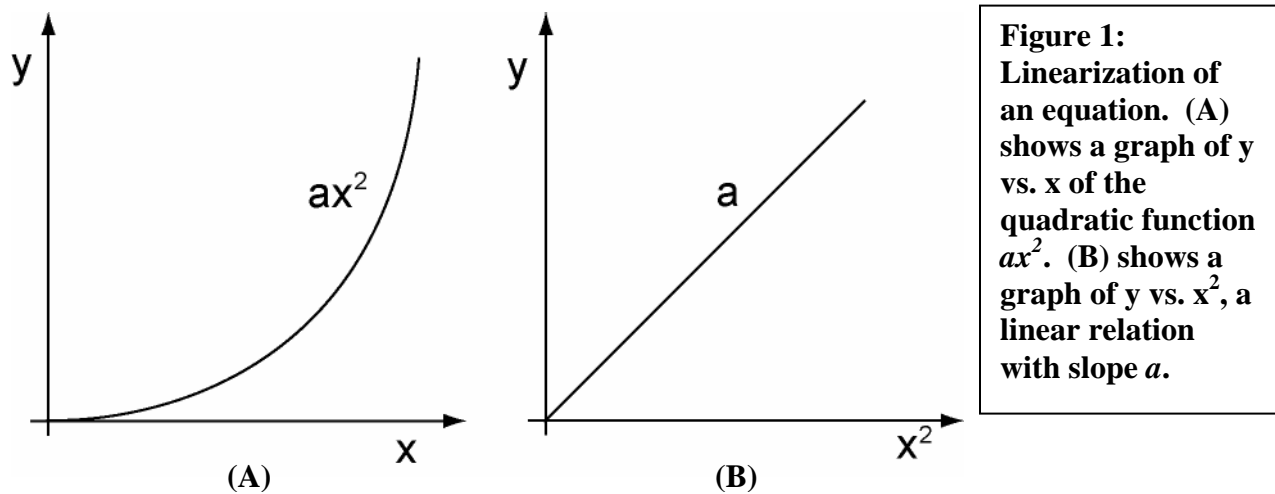
What's the reason for wanting to do a Least Squares Fit? Why bother with finding a best fit line to a set of data in the first place? Well, a graph is used to show whether there is a relation between the dependent variable (the y-axis) and the independent variable (the x-axis). Usually, one looks for a linear relation, that is to say whether the data points fall roughly on a line or not. A linear relation has the form  $y = a + bx$ , which is useful for showing direct relationships such as  $F=ma$  and  $V=IR$ . A graph of the force of gravity vs. mass would yield a line with a slope equal to the acceleration due to gravity. A graph of voltage vs. current would give a value for resistance. This is good stuff!

What about equations which are non-linear? How could calculating a best fit line using the Least Squares Fitting method help with that? Here are two examples of equations that may appear non-linear but can be made linear.

The first example involves the magnetic force on electrons and the circular motion the electrons undergo in a uniform magnetic field. The equation is  $eB = \frac{m}{r} \sqrt{\frac{2eV}{m}}$ , which doesn't seem like it could be graphed easily. However, one wants to graph the charge vs. the mass of the electron, as was done in 1897 by Sir J. J. Thomson. Therefore, that messy equation can be rearranged as follows:

$$\frac{eBr}{m} = \sqrt{\frac{2eV}{m}} \longrightarrow \frac{e^2 B^2 r^2}{m^2} = \frac{2eV}{m} \longrightarrow \frac{eB^2 r^2}{m} = 2V \longrightarrow e = \frac{2V}{B^2 r^2} m$$

This involves simple algebra, however, one had to know ahead of time what variables one wanted to graph. In order to graph charge vs. mass, the charge had to be alone on the left-hand side, and the mass had to be on the right-hand side, to the first power only, and with a prefactor of known variables ( $V$ ,  $B$ , and  $r$ ). This gives a linear graph of the charge vs. the mass with a slope of  $\frac{2V}{B^2 r^2}$ .



The second example involves the period of a pendulum, given by  $T = 2\pi\sqrt{l/g}$ . In lab one measures the period and the length of the pendulum. What to do? Squaring the equation gives  $T^2 = \frac{4\pi^2}{g}l$ .

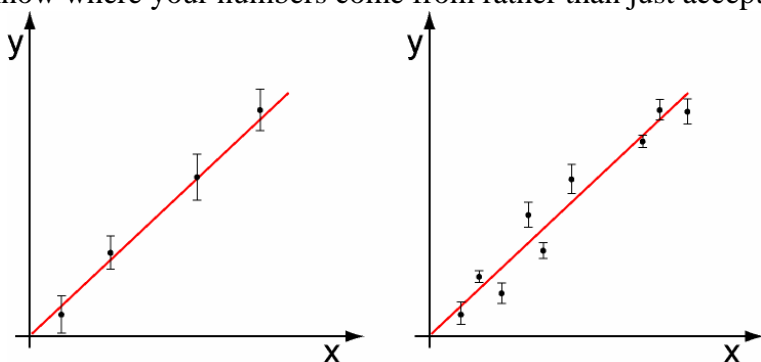
This yields a linear graph of the square of the period vs. the length of the pendulum with a slope of  $\frac{4\pi^2}{g}$ .

## 2. Methods of Finding the Best Fit Line: Estimating, Using Excel, and Calculating Analytically

A graph can be used to estimate the best fit line as opposed to calculating the best fit line from the data points. The basic idea is to draw a line through the data with as many data points above it as below it within error. For full instructions refer to the **Making Graphs** section in the lab manual. Estimating the best fit line is a good idea when there are fewer than 5 data points. Picking where the line should go is a skill that improves with practice. The downside of this is that the best fit line and the uncertainties really are just estimates based on “eye-balling” the graph.

A step up from getting a best fit line by hand is having Excel graph and fit a **trendline** (mentioned in the **Graphing with Excel** document). Excel uses the Least Squares Fit method to calculate the best fit line. Using Excel is a good idea for data sets larger than 5 points, for the program takes care of the whole process; but, while an Excel graph will give the equation of the best fit line, it won't give the uncertainty in the slope. Also, it is always dangerous to use a result that is not fully understood.

Calculating the best fit line by using Least Squares Fit method is good for data sets with more than 5 points and gives better results for more data points. The analytic method is good for when error is small. The method will give numbers for slope, intercept, the uncertainties in the slope and intercept, as well as the correlation coefficient (an indication of how good the data fit the line). This document goes through the derivation and the method of doing a least squares fit. At the end, there are instructions for using Excel's LINEST function, which calculates all of the numbers of the fit for you. Why bother reading this whole document when you could just skip to the end? It's important to know where your numbers come from rather than just accepting them from a program.



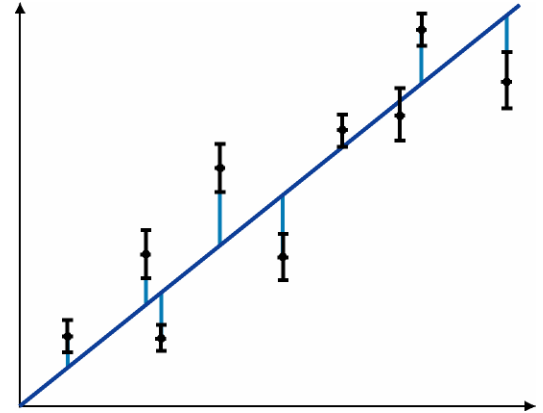
**Figure 2: Comparison of (A) data where it's easy to estimate the best fit line and (B) data where it's best to use a Least Squares Fit.**

The key to calculating a Least Squares Fit is a well-organized data sheet (tired of hearing that?). Equations will be used to calculate the values for the slope, the intercept, and the uncertainty in those values. Throughout the document, the equations will refer to the independent variable as just  $x$  and the dependent variable as just  $y$ . When applying Least Squares Fitting to a data set, keep in mind

which variable is the *independent* one (the variable one changes in lab) and which one is the *dependent* one (the variable one checks in lab to see what effects changes in the independent variable produce).

### 3. How to Calculate a Least Squares Fit

Consider the distances from each point to the best fit line, as shown in **Figure 3**, also called the deviations. If a line is a really good fit, those deviations will be as small as possible. The Least Squares Fitting method attempts to minimize the square of the deviations. (Squaring the deviations makes the math far more manageable, where working with the absolute values of the distances would lead to discontinuous derivatives.) The sum of all of the squares of the deviations is called the residual,  $\chi^2$  ( $\chi$  is pronounced “ky” to rhyme with “guy”), given by equation (1), where  $y_i$  are the data points and  $y_{true}$  is the  $y$ -value from the best fit line.



**Figure 3: Example of deviations from the best fit line (blue).**

$$\chi^2 \equiv \sum_{i=1}^N (y_i - y_{true})^2 \quad (1)$$

If the data follows a linear relation, then  $y_{true}$  can be expressed in the general form  $y_{true} = mx + b$ , making the residual

$$\chi^2 = \sum_{i=1}^N (y_i - (mx_i + b))^2 = \sum_{i=1}^N (y_i - mx_i - b)^2 \quad (2)$$

Note that this formula still uses the general form of the equation for the line:  $b$  and  $m$  have not been specified! To find the line that best fits the data, the residual should be as small as possible.

Variables  $b$  and  $m$  must be chosen so that they minimize  $\chi^2$ . This is where the method gets its name: making the sum of the squares of the deviations the least it can be. Any statistical book will give the details of this minimization; they are also available on the website under **Useful Documents and Websites**. The end result is that the values of  $b$  and  $m$  are given by the following equations:

$$b = \frac{\left( \sum_{i=1}^N (x_i^2) \right) \left( \sum_{i=1}^N (y_i) \right) - \left( \sum_{i=1}^N (x_i) \right) \left( \sum_{i=1}^N (x_i y_i) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2} \quad (3)$$

$$m = \frac{N \left( \sum_{i=1}^N (x_i y_i) \right) - \left( \sum_{i=1}^N (x_i) \right) \left( \sum_{i=1}^N (y_i) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2} \quad (4)$$

#### 4. Uncertainty in the Dependent Variable, Slope, and Intercept

The formula for the standard deviation is  $\sigma_x = \sqrt{\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2}$ . That formula applies to finding the standard deviation of a number of measurements of the same “true” value, which are randomly distributed around that “true” value (assuming that systematic errors have been reduced). An example of that would be finding the distance to a block sitting on the table: various measurements of the distance might be slightly off the “true” distance to the block, but that “true” distance is the same for all the measurements because the block is sitting still. The “true” value for the distance is approximated by the average of the distance measurements.

Consider the measurement of the distance to a glider as it moves along a track. Time is the independent variable ( $x$ ) and the distance to the glider is the dependent variable ( $y$ ). In this example these data are separate measurements in themselves, each one randomly distributed around the “true” value for the distance at a particular instant. If lots of measurements of the distance to the glider were taken *at one instant* by lots of rangers instead of just one, those measurements would be randomly distributed around a “true” value for the distance to the glider *at one instant in time*. The “true” value for the distance to the glider at that instant could be approximated by the average value of the measurements from all those rangers --- but that’s just one instant!

The mean of the distance ( $\bar{y}$ ) is a good approximation to the “true” value of the distance *if the glider is not moving*. If the glider was given a push and is moving along a track, measurements of distance versus time would follow a linear pattern. As seen in **Section 3**, the best fit line ( $y_{true} = mx + b$ ) is a good approximation to the “true” value of the distance at any time *if the glider is moving*.

The standard deviation of the dependent variable,  $\sigma_y$ , is calculated a little differently than  $\sigma_x$ . When the dependent variable is changing with the independent variable in a linear fashion, the standard deviation of the dependent variable can be found by the following:

$$\sigma_y = \sqrt{\frac{1}{N-2} \sum_{i=1}^N (y_i - y_{true})^2} = \sqrt{\frac{1}{N-2} \sum_{i=1}^N (y_i - mx - b)^2} \quad (5)$$

In the formula for the standard deviation,  $\sigma_x$ , the factor before the sum is  $\frac{1}{N-1}$ , but in equation (5)

for  $\sigma_y$  it is  $\frac{1}{N-2}$ . Why the difference? Consider a data set with only two points (i.e.  $N = 2$ ). Two points define a line, so (of course!) the line through those two points is a perfect fit. For a data set consisting of less than three points, the uncertainty  $\sigma_y$  is undefined.

The uncertainty in  $b$ , the intercept, and  $m$ , the slope, can be found by propagation of the uncertainty in the dependent variable,  $\sigma_y$ , in the formula for  $b$  and  $m$ . In doing this, it’s assumed that most of the error comes from the dependent variable. This propagation may be found in any statistics book and is also available on the website under **Useful Documents and Websites**.

The uncertainties are given by the following equations:

$$\sigma_b = \sqrt{\frac{\sigma_y^2 \left( \sum_{i=1}^N (x_i^2) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}} \quad (6)$$

$$\sigma_m = \sqrt{\frac{N \sigma_y^2}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}} \quad (7)$$

## 5. Calculating the $r^2$ Value

Thanks to the efforts of everyone from Mr. Babbage to Mr. Gates, it's no longer necessary to watch the best years of your life slip past while you hand-process large data sets. In fact, at the push of a button we can even see how well any TWO data sets relate to each other by linear regression. But as you might suspect, the quality of all this cyber-math must be reported with something more objective than the words "good" or "bad". So when your linear regression program spits out a slope and a y-intercept, it probably also gives you a number labeled either "r" or "correlation coefficient". What is that number?

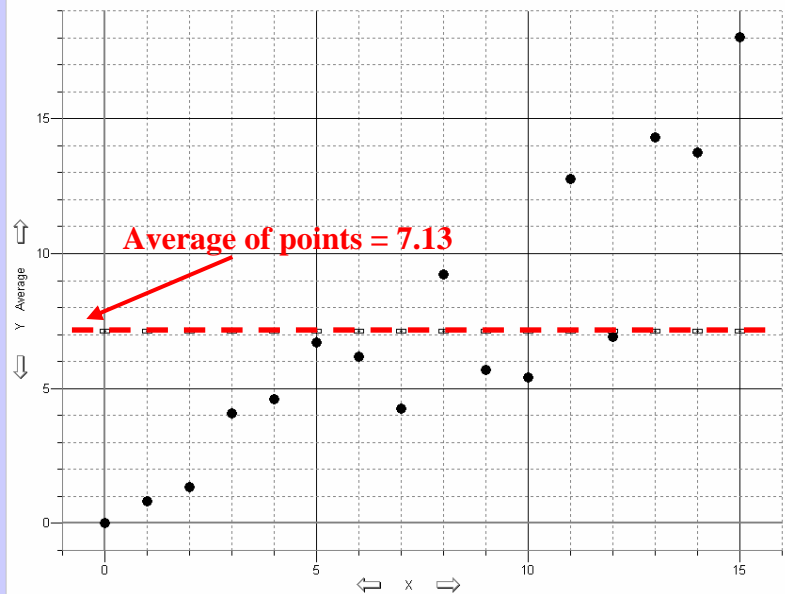
The first thing you need to know about it is that it ONLY tells you how well the two data sets match a STRAIGHT LINE. It does not tell you how well the data matches any other mathematical function.

The second important thing to know is that the correlation coefficient has a range from 0 to 1. If every matched pair of data values is a point that fits exactly on a single line, r will equal 1; a perfect fit. If the points are close to the line but not perfect, r will be less than 1 by a small amount. Finally, if the points are all over the place, and not even close to favoring the chosen line, r will equal 0. This means there is no LINEAR relation between the two data sets.

A third thing to know is that we've been using the coefficient's nickname. Its full name is the "Pearson product-moment correlation coefficient". Knowing this may get you past a tough question on Jeopardy some day.

Now let's go back to that second thing. How do you get a calculation to be 1 when the points match the line and 0 when they miss?

Start by calculating the average of all the "y" values ( $\bar{y}$ ). On a graph, this "y" average will be a horizontal line that runs through all the data points around the mid-height point, as shown by the dotted line in **Figure 4**. Why start with the average? The horizontal average line is no better



**Figure 4: Average of dependent variable values.**

than any other guess as a model for the line that fits the data. It gives a “worst case scenario” to use as a basis of comparison. To see how good (or bad) this fit is, you could take the difference between each actual  $y$  value ( $y_i$ ) and the average ( $\bar{y}$ ), and add them up...

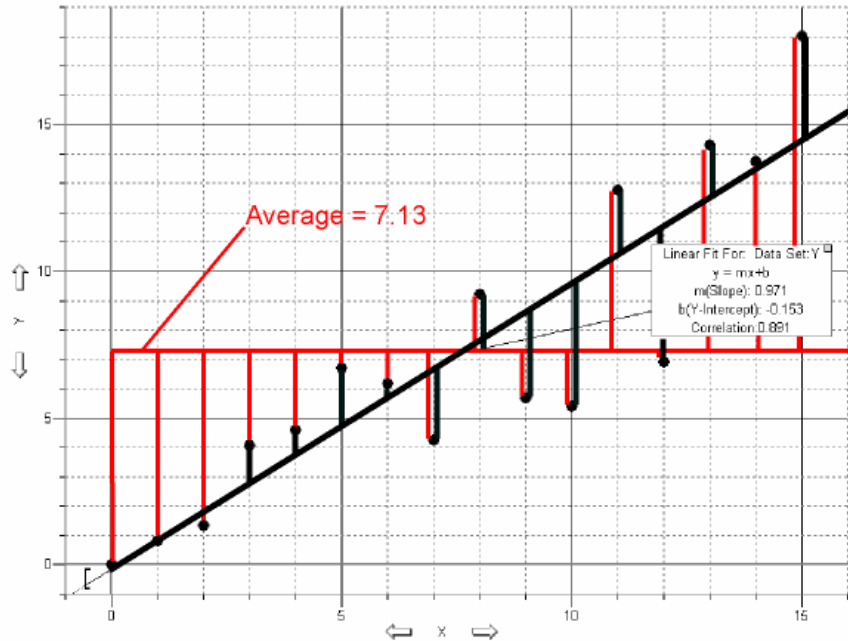
Oops! That’s always going to be zero by definition. The average is the value that’s right in the middle of a data set. Half the values are larger and half are smaller, so the sum of positive differences will exactly match the negative differences. To remedy this, after finding each difference, square it (to make all values positive) and then add them all up (since we are only using this number for comparison, we’ll just remember to square all other values that will be compared to it):

$$\sum (y_i - \bar{y})^2 \quad (8)$$

This sum of the square of the differences is a good way of determining how good the fit is. A small value for this sum indicates all the numbers are all pretty close to the average. A large value indicates the numbers are widely scattered. This sum will be used as a comparison for any additional attempts to find a line that fits the data.

You then calculate a “ $y_{calc}$ ” value for each  $x$  value using the equation of the line you wish to fit to the data:

$$y_{calc} = mx + b$$



**Figure 5: Straight line fit (black) to data with the equation of the fit  $y = 0.971x - 0.153$  and correlation coefficient 0.891. Average of **y-values** (red) is 7.13. Note the smaller deviations for the best fit line.**

To see how well this line matches the data, a similar sum of squared differences is calculated.

$$\sum (y_i - y_{calc})^2 \quad (9)$$

It might seem that a simple comparison with equation (8) might be the way to evaluate how good this fit is, but we can do much better. Consider subtracting equation (9) from equation (8) to get equation (10).

$$\sum (y_i - \bar{y})^2 - \sum (y_i - y_{calc})^2 \quad (10)$$

If the chosen line  $y_{calc} = mx + b$  fits no better than the horizontal average line of the numbers, then equation (9) will be just as large as equation (8). If that is the case, the result of (10) will be zero.

$$\sum (y_i - \bar{y})^2 - \sum (y_i - y_{calc})^2 = 0 \quad (11)$$

If the chosen line  $y_{calc} = mx + b$  is an exact fit, then each calculated value ( $y_{calc}$ ) is exactly equal to the corresponding data value ( $y_i$ ), and equation (9) will be zero. If that is the case, the result of equation (10) will be equation (8).

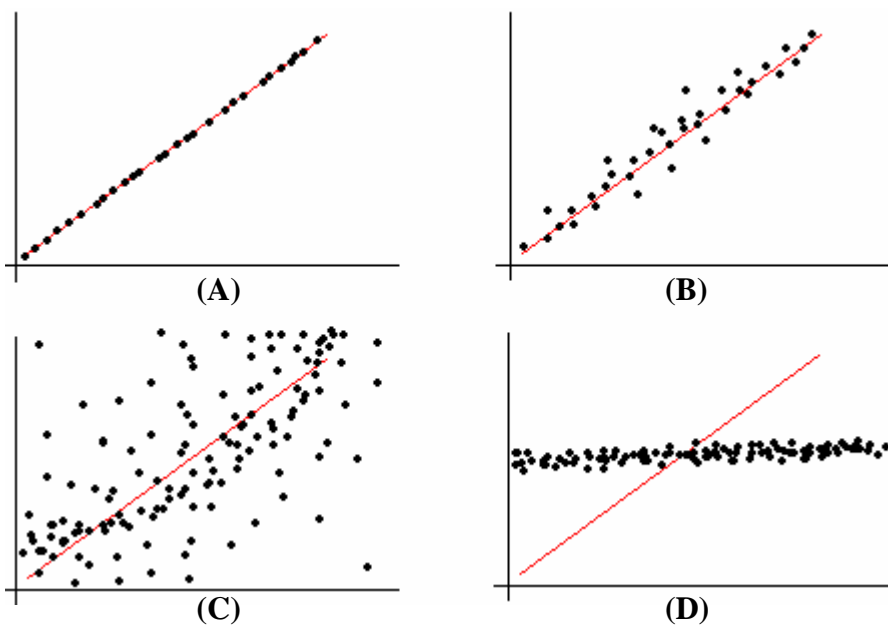
$$\sum (y_i - \bar{y})^2 - 0 = \sum (y_i - \bar{y})^2 \quad (12)$$

So equation (10) is equal to zero for a line that completely misses the data and equal to equation (8) for an exact match to the data. If we divide equation (10) by equation (8), then it will equal zero for a line that completely misses and equal 1 for an exact match. This process of dividing equation (10) by equation (8) is called normalization. When normalized, the correlation coefficient,  $r^2$ , is adjusted to have a range from 0 to 1.

$$r^2 = \frac{\sum (y_i - \bar{y})^2 - \sum (y_i - y_{calc})^2}{\sum (y_i - \bar{y})^2} \quad (13)$$

You may see other formulae for  $r$  that involve standard deviations. They are alternate methods to calculate the same number. In the past, calculators had only standard deviation buttons so these other formulae may have been easier to use. Today it is common for the correlation coefficient to be built into the software, so the difficulties of calculation are no longer an issue. We have developed the equation using averages because it is much easier to see why it takes the form it does.

**Figure 6A** is almost a perfect fit, for the points fall very close to the best fit line. **Figure 6B** is a good fit, for the points fall fairly close to the best fit line. **Figure 6C** is a bad fit, for the points do not fall close to the best fit line (and if you look closely you'll note that they seem to fall in a



**Figure 6: Graphs with different  $r^2$  values. (A) is a perfect fit. (B) is a good fit. (C) is a bad fit, not linearly related. (D) is a bad fit, zero relationship between the variables.**



parabola and not a line!). **Figure 6D** is a bad fit for the points do not fall along the chosen best fit line. The data in **Figure 6D** falls on a horizontal line, which indicates that the variables are unrelated, for changes in the independent variable produce no changes in the dependent variable.


## 6. Sample Data and Setting up the Spreadsheet


Here is a run-through of the method of Least Squares Fitting using the sample data shown in **Figure 7** according to the tables in **Section 6.3**. Notice that the first row is **boldfaced**, as well as rows 13 through 23. Titles of columns should always be **boldfaced**, and it helps the results and their titles stand out from the data.

	A	B	C	D	E	F	G	H	I	J	K
1	<b>x</b>	<b>y</b>	<b>x<sup>2</sup></b>	<b>y<sup>2</sup></b>	<b>xy</b>	<b>dev x</b>	<b>dev y</b>	<b>(dev x)<sup>2</sup></b>	<b>(dev y)<sup>2</sup></b>	<b>(dev x)(dev y)</b>	<b>(y-b-mx)<sup>2</sup></b>
2	1	4	1	16	4	-5	-24	20	557	106	0
3	2	9	4	72	17	-4	-19	12	365	67	1
4	3	17	9	289	51	-3	-11	6	112	27	6
5	4	19	16	342	74	-2	-9	2	83	14	2
6	5	22	25	484	110	-1	-6	0	31	3	9
7	6	36	36	1296	216	1	8	0	71	4	34
8	7	32	49	1024	224	2	4	2	19	7	11
9	8	43	64	1849	344	3	15	6	237	39	6
10	9	44	81	1936	396	4	16	12	269	57	3
11	10	51	100	2601	510	5	23	20	548	105	0
12											
13			<b>sum x</b>	<b>55</b>		<b>N</b>	<b>10</b>				
14			<b>sum y</b>	<b>276</b>		<b>mean x</b>	<b>5.5</b>				
15			<b>sum(x<sup>2</sup>)</b>	<b>385</b>		<b>mean y</b>	<b>27.6</b>				
16			<b>sum(y<sup>2</sup>)</b>	<b>9910</b>							
17			<b>sum(xy)</b>	<b>1946</b>		<b>b (intercept) =</b>	<b>-0.93</b>				
18			<b>sum(dev x)</b>	<b>0</b>		<b>m (slope) =</b>	<b>5.19</b>				
19			<b>sum(dev y)</b>	<b>0</b>		<b>σ<sub>y</sub> =</b>	<b>2.99</b>				
20			<b>sum[(dev x)<sup>2</sup>]</b>	<b>83</b>		<b>σ<sub>b</sub> =</b>	<b>2.04</b>				
21			<b>sum[(dev y)<sup>2</sup>]</b>	<b>2292</b>		<b>σ<sub>m</sub> =</b>	<b>0.33</b>				
22			<b>sum[(dev x)(dev y)]</b>	<b>428</b>		<b>r<sup>2</sup> =</b>	<b>0.97</b>				
23			<b>sum[(y-b-mx)<sup>2</sup>]</b>	<b>71</b>							

**Figure 7: Screen capture of the worksheet constructed by following the steps in this document. Note that the means have been formatted to show only one decimal place, the best fit line calculations only two decimal places, and the remaining numbers no decimal places. This was done to make the worksheet easier to read.**

### 6.1 Example: Entering a Sum

Cell D13 should have the sum of (A2:A11). Any time a formula is entered in Excel, precede it by an equal sign; otherwise, Excel thinks the entry is just text. Click on D13, type an equal sign in the cell, go to the **Name Box**, and select the **SUM** function. In the **Function Arguments** window that appears, use the **Number 1** field to select all of the independent variable data. Do this by clicking the  button next to the field. The **Function Arguments** window will collapse, leaving only the **Number 1** field showing. Use the mouse pointer to select all of the independent variable data,

meaning cells A2 through A11. A dashed box indicates what data is selected. If a mistake is made, click on the **Number 1** field, delete the information, and try selecting again. When done selecting data, click the  button next to the **Number 1** field, which will expand the **Function Arguments** window. Click the **Ok** button.

The short method of doing this is just to type =SUM(A2:A11).

## 6.2 Example: Entering a Formula

Finding the sum of the square of  $x_i$ ,  $\sum_{i=1}^N (x_i^2)$ , requires first setting up a column C which calculates the squares. In cell C2 the formula for  $(A2)^2$  is A2^2. To fill column C from C2 all the way to C11 with the squares of the respective cells (A2:A11), drag the **fill handle** (mentioned in the **Introduction to Excel** document, **Section 2.4 Filling Data**). Give the column the title of  $x^2$  in cell C1.

## 6.3 Formulae for Excel Using Sample Data

Set up the columns A through K as described in **Figure 8**. These columns will be used to calculate the necessary sums.

<b>Figure 8: Formulae for the columns in Excel using sample data from Figure 7.</b>		
<b>Column</b>	<b>Rows Contain</b>	<b>Sample Excel Formula</b>
A	$x_i$	n/a
B	$y_i$	n/a
C	$x_i^2$	A2^2
D	$y_i^2$	B2^2
E	$x_i y_i$	A2*B2
F	$x_i - \bar{x}$	A2-\$G\$14
G	$y_i - \bar{y}$	B2-\$G\$15
H	$(x_i - \bar{x})^2$	F2^2
I	$(y_i - \bar{y})^2$	G2^2
J	$(x_i - \bar{x})(y_i - \bar{y})$	F2*G2
K	$(y_i - b - mx_i)^2$	(B2-\$G\$17-\$G\$18*A2)^2

Calculate the sums in the cells using **Figure 9** using the sample data from **Figure 7**. These sums will be used to calculate the slope, the intercept, the uncertainties, and the correlation coefficient.

Calculate the values for the Least Squares Fit according to the formulae in **Figure 10**.

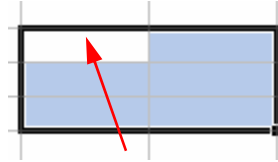
<b>Figure 9: Formulae for the cells which calculate the sums.</b>	
Cell	Formulae
D13	$\sum_{i=1}^N (x_i)$ SUM(A2:A11)
D14	$\sum_{i=1}^N (y_i)$ SUM(B2:B11)
D15	$\sum_{i=1}^N (x_i^2)$ SUM(C2:C11)
D16	$\sum_{i=1}^N (y_i^2)$ SUM(D2:D11)
D17	$\sum_{i=1}^N (x_i y_i)$ SUM(E2:E11)
D18	$\sum_{i=1}^N (x_i - \bar{x})$ SUM(F2:F11)
D19	$\sum_{i=1}^N (y_i - \bar{y})$ SUM(G2:G11)
D20	$\sum_{i=1}^N (x_i - \bar{x})^2$ SUM(H2:H11)
D21	$\sum_{i=1}^N (y_i - \bar{y})^2$ SUM(I2:I11)
D22	$\sum_{i=1}^N ((x_i - \bar{x})(y_i - \bar{y}))$ SUM(J2:J11)
D23	$\sum_{i=1}^N ((y_i - b - mx)^2)$ SUM(K2:K11)
G13	$N$ n/a
G14	$\sum_{i=1}^N \left( \frac{x_i}{N} \right)$ D13/G13
G15	$\sum_{i=1}^N \left( \frac{y_i}{N} \right)$ D14/G13

<b>Figure 10: Formulae for the cells which calculate the <math>mx_i+b</math>, <math>\sigma_y</math>, <math>\sigma_b</math>, <math>\sigma_m</math>, and <math>r^2</math> value.</b>	
Cell	Formulae
G17	$b = \frac{\left( \sum_{i=1}^N (x_i) \right) \left( \sum_{i=1}^N (y_i) \right) - \left( \sum_{i=1}^N (x_i y_i) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}$ (D15*D14-D13*D17)/(G13*D15-D13*D13)
G18	$m = \frac{N \left( \sum_{i=1}^N (x_i y_i) \right) - \left( \sum_{i=1}^N (x_i) \right) \left( \sum_{i=1}^N (y_i) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}$ (G13*D17-D13*D14)/(G13*D15-D13*D13)
G19	$\sigma_y = \sqrt{\frac{1}{N-2} \sum_{i=1}^N (y_i - b - mx_i)^2}$ SQRT((1/(G13-2))*D23)
G20	$\sigma_b = \sqrt{\frac{\sigma_y^2 \left( \sum_{i=1}^N (x_i^2) \right)}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}}$ SQRT((G19^2*D15)/(G13*D15-D13*D13))
G21	$\sigma_m = \sqrt{\frac{N \sigma_y^2}{N \left( \sum_{i=1}^N (x_i^2) \right) - \left( \sum_{i=1}^N (x_i) \right)^2}}$ SQRT((G13*G19^2)/(G13*D15-D13^2))
G22	$r^2 = \frac{\sum (y_i - \bar{y})^2 - \sum (y_i - y_{calc})^2}{\sum (y_i - \bar{y})^2}$ (D21-D23)/D21

## 6.4 Using Excel's LINEST Function

Congratulations for making it through this exercise in Least Squares Fitting! This exercise set up quite a spreadsheet. There is, in fact, a much shorter way of getting the values of the Least Squares Fit, the uncertainties, and the correlation coefficient, involving an Excel function. Why make the spreadsheet if there is a function that already does the fit? Going through the work of setting up the spreadsheet creates a familiarity with the method of Least Squares Fitting. As was mentioned before, it's dangerous to use a result without understanding it.

The LINEST function can return a single value for the slope of the line or it can return all of the values of the Least Squares Fit. It takes as arguments  $\{y_n\}$  (the dependent variable data),  $\{x_n\}$  (the independent variable data), an indicator which can set the intercept of the equation to zero, and an indicator which can display the values of the fit in addition to the slope. Before using the function, enter the independent and dependent variable data in columns A and B, just like in **Figure 7**. Now, select six blank cells in a block, two columns wide and three rows high. Make sure the upper left-most cell is the first cell clicked when making the selection, as shown in **Figure 11**.



**Figure 11: Selecting the cells for the LINEST function. Red arrows points to the cell which was clicked first and where the formula will be entered.**

Click **F2** and enter the formula `LINEST(B2:B11,A2:A11,,TRUE)`. When done entering the formula, press the **CTRL** key, **SHIFT** key, and **ENTER** key all at once; this will tell Excel to write the results of the LINEST function in the selected cells with the pattern shown in **Figure 12**. Note that Excel will only display the numerical results; keep in mind which cell refers to which quantity!

<b>m</b>	<b>b</b>	5.188	-0.9333
$\sigma_m$	$\sigma_b$	0.329	2.04209
$r^2$	$\sigma_y$	0.969	2.98931

(A)

(B)

**Figure 12: (A) Pattern of results from LINEST function, (B) what Excel will display.**