# Data Analysis with Spreadsheets

## 1 INTRODUCTION AND LEARNING GOALS

Spreadsheets are used throughout this course as the primary tool to record, interpret, and present the measurements that you make in lab.

At the end of this assignment, you should understand:

- Basic spreadsheet functionality including the use of custom and built-in functions.

- How to make plots that meet the style requirements for this course.

A spreadsheet application is required to complete the assignment. Instructions for the assignment are written for Microsoft Excel. Johns Hopkins students have free access to Microsoft Excel at jhu.onthehub.com. Other spreadsheet programs may be sufficient to complete the activity, but students should be aware that use of programs besides Excel will be at your own risk. Lab TAs and staff will not be able to provide support and technical assistance for programs other than Excel.

## 2 BACKGROUND

### 2.1 READING ASSIGNMENT

Read Chapters 2.1-2.6 (pp. 13-28) in Taylor, J. R. (1997) *An Introduction to Error Analysis.* Sausalito, CA: University Science Books.

### 2.2 FITTING A LINEAR MODEL TO DATA

Once data have been plotted it is very common to interpret the observations using a mathematical model. One can do this by:

1. Adding a mathematical model to the plot of the data (see Figures 2.2 and 2.3).

2. Optimize the parameters of the mathematical model to obtain a the best possible fit to the data.

A simple example of this procedure is fitting a straight line model to data.

The mathematical model for a line is given by:

$$y = ax + b \tag{2.1}$$

where $a$ is a parameter that describes the slope and $b$ is the $y$-intercept of the line.

In Excel one can build a model to fit to data by setting line parameters in cells and creating an additional column with the predicted line based on those cells. For example, looking at Figure 2.1, one might guess that the slope is approximately:

$$a_{guess} = \frac{-5}{40} = -0.125. \tag{2.2}$$

Similarly, a reasonable guess at the $y$-intercept might be:
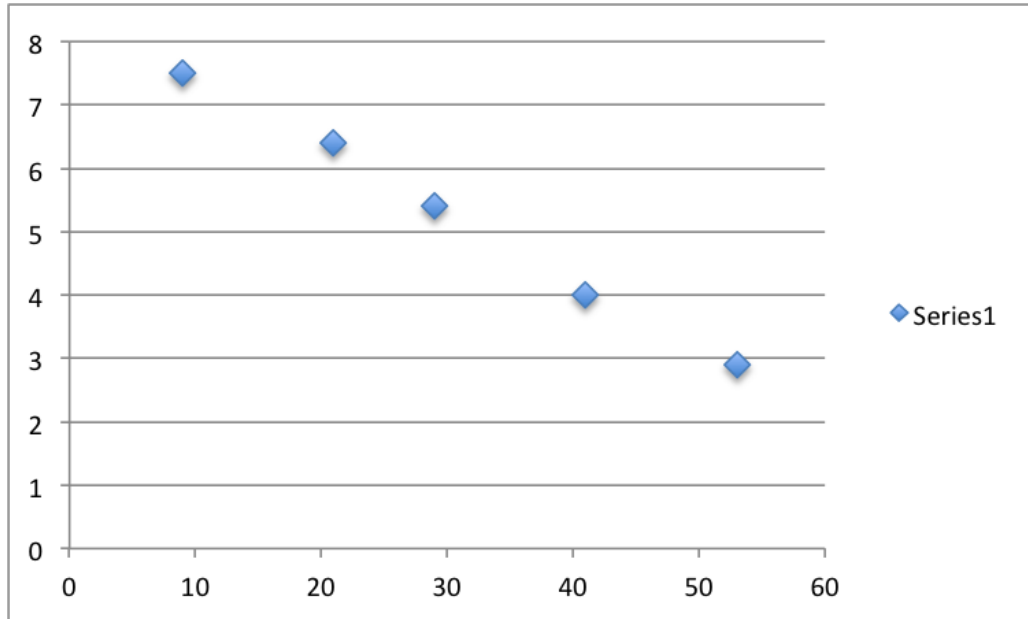
$$b_{guess} = 9. \tag{2.3}$$



Figure 2.1: A crude first pass at plotting the data. This plot is poorly labeled, difficult to read, and does not describe the data it is meant to convey.

One could then construct the line as shown in Column G of Figure 2.2. This "guessed" line can be plotted on the data as shown in Figure 2.3. We see that the original guess at the slope and intercept parameters were not too far off; the model appears to roughly agree with the data. However, we also notice that our model seems to over-estimate, when compared to the data, at low values of $x$ and under-estimate at high values.

To optimize the model, a variable can be constructed to quantify how well the model matches the data. One approach is to measure the vertical discrepancy between the prediction of the model and the observed value for each data point. For example:

$$\chi^2 = \sum \left( y_{i\,\text{data}} - y_{i\,\text{model}} \right)^2. \tag{2.4}$$

where $y_{i\,\text{data}}$ are the observed values in data and $y_{i\,\text{model}}$ are the predictions based on the chosen model. In the case of the current example, the model is a straight line.

When the model and the data agree perfectly, Equation 2.4 will be zero. The better the model matches the data, the smaller $\chi^2$ will be. So, by continuously adjusting the slope and intercept parameters, we can find a line of best fit. Since we are minimizing $\chi^2$, this process is often called a "least-squares fit".

Excel can vary the parameters of the model for you to perform a least-squares fit. The process is as follows:

1. Calculate $\chi^2$ using an *array formula*. Using the cells shown in Figure 2.2, Equation 2.4 can be written in spreadsheet syntax as:

    `=SUM((D2:D6-G2:G6)^2)`

    Instead of simply typing ⏎, one needs to tell Excel to make this calculation as an *array formula*. That is, we want Excel to calculate:

    `(D2-G2)^2 + (D3-G3)^2 + (D4-G4)^2 + ...`

    Execute the formula as an *array formula* by pressing ⎈Ctrl + ⇧ + ⏎ (Windows) or ⌘ + ⇧ + ⏎ (Mac). Note that ⇧ denotes the "shift" key *not* the "up arrow" key. Note that when an array formula is successfully executed, Excel will automatically enclose the entered equation in braces.

2. Load the Solver Add-in in Excel. Follow the instructions linked below.

    - For Windows: "Load the Solver Add-in"

    - For Mac: "Load the Solver Add-in in Excel 2016 for Mac".

3. Use Solver to Minimize $\chi^2$. Excel offers a tool for automatically varying cell values. Depending on your platform, Solver can be found in the following places:

    - Windows: The Solver command is available in the Analysis group on the Data tab.

    - Mac: Use the menu: Tools 〉 Solver...

    The Solver tool is configured with the box shown in Figure 2.4.

    **Set Objective:** The "Objective" is the target cell that Solver will attempt to minimize (set by the "Min" option). In this example, the Objective is set to cell H2.

    **by Changing Variable Cells:** The variables that Solver should vary are specified. When more than one cell should be varied, cells are specified in a comma-separated list. In our example, the "Variable Cells" are set to Sheet1!$F$3 and Sheet1!$F$5.

    **Make Unconstrained Variables Non-Negative:** In most cases, we will not want to constrain any variables to be non-negative. We will make sure that this box is not checked.

    **Solving Method:** The "Solving Method" gives users the option to select which computational algorithm is used in the calculation. The default, "GRG Nonlinear" option typically works well.

If Solver is able to converge and successfully find a minimum, a dialog box, like that shown in Figure 2.5) will report the result. By clicking "OK" the solution that Solver found will appear in the variable cells.

For more detailed help using Solver, see Microsoft's help article: "Define and solve a problem by using Solver".

By performing the least-squares fit – with the steps described above – it was found that the slope and intercept parameters in this example are:

$$a = -0.107172125 \simeq -0.11 \tag{2.5}$$

$$b = 8.519467109 \simeq 8.5. \tag{2.6}$$

After running Solver, the line that was originally plotted as the "guessed" model will be updated with the new parameters. The updated fit line is shown in Figure 2.6. It is clear that the fit line does a much better job modeling the data than the original crude guess. Note also that the fit parameters $a$ and $b$, match those that are returned by LINEST() shown in Figure A.4

| G2 | | | $f_x$ | =$F$3*B2+$F$5 | | | |
|---|---|---|---|---|---|---|
| | A | B | C | D | E | F | G |
| 1 | Trial (unitless) | Xerbert (a) | Xerb. Error (a) | Yeehaw (b) | Yeehaw Error (b) | Fit Parameters | Fit Line |
| 2 | 1 | 9 | 2 | 7.5 | 0.7 | Slope | 7.875 |
| 3 | 2 | 21 | 2 | 6.4 | 0.6 | -0.125 | 6.375 |
| 4 | 3 | 29 | 2 | 5.4 | 0.3 | Intercept | 5.375 |
| 5 | 4 | 41 | 2 | 4 | 0.9 | 9 | 3.875 |
| 6 | 5 | 53 | 2 | 2.9 | 0.5 | | 2.375 |
| 7 | | | | | | | |

Figure 2.2: A linear model can be constructed in Excel by setting the slope and intercept parameters (cells F3 and F5 respectively). The model is then computed by calculating the resulting $y$-value with the *observed* $x$-values. Note that the equation anchors the slope and intercept cells with dollar-signs ("$F$3" and "$F$5"). The dollar signs tell Excel to never change these cells when this formula is applied to adjacent cells.
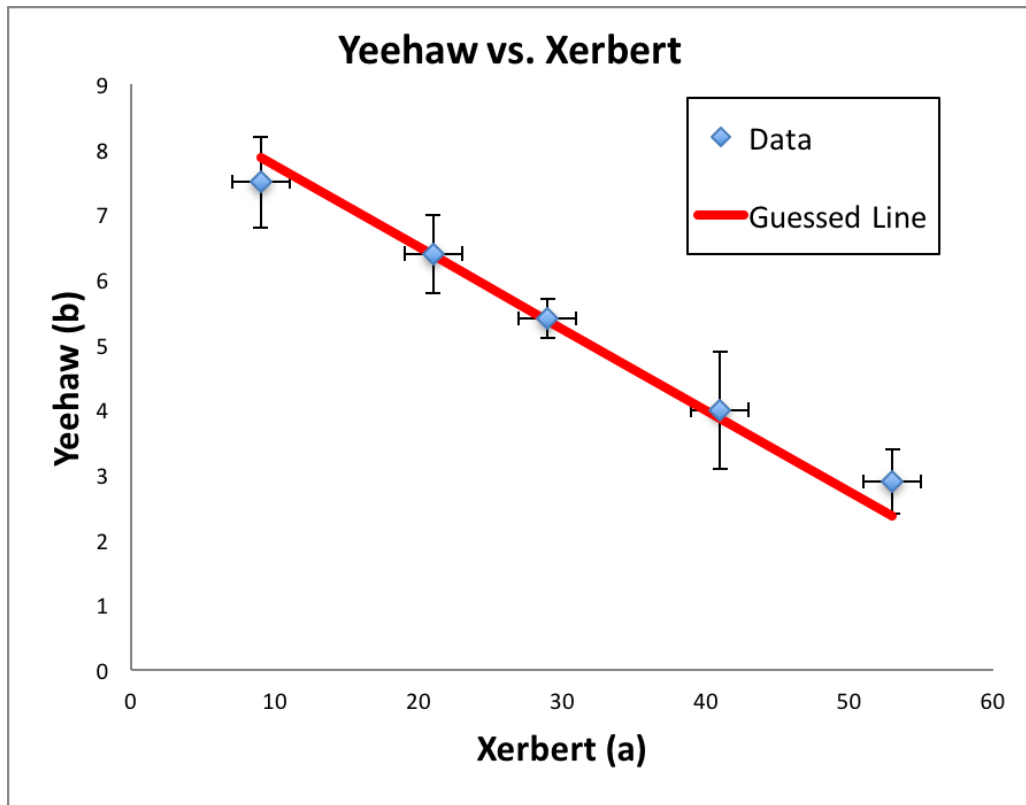
Figure 2.3: An attempt to model the observed data with a straight line. The slope and intercept parameters are just guesses. The central three data points appear to match the model very well. At low and high values of $x$ though, the model appears to do a worse job.
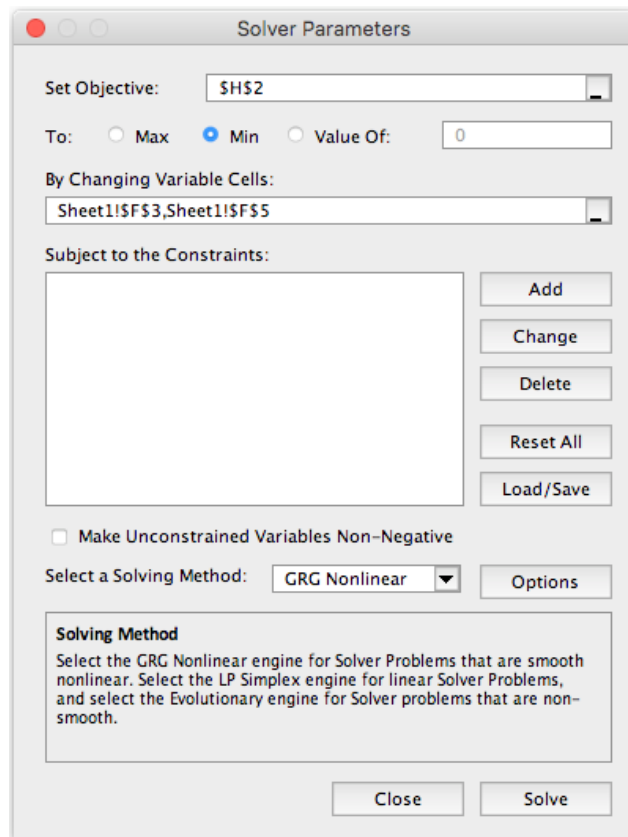
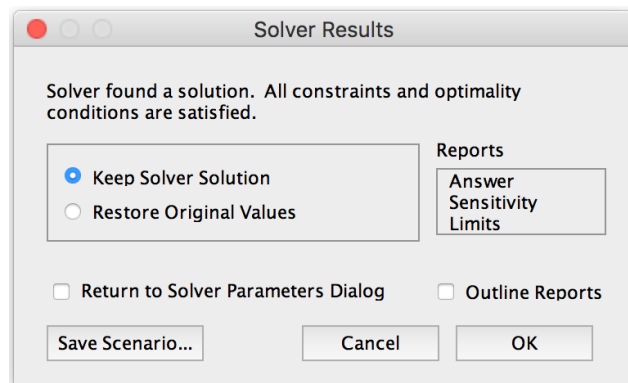Figure 2.4: The Solver Parameters dialog box is launched when Solver is run.



Figure 2.5: When Solver is finished, an informational dialog box will report success or failure of the selected operation. Clicking "OK" will keep the solver solution that was found.
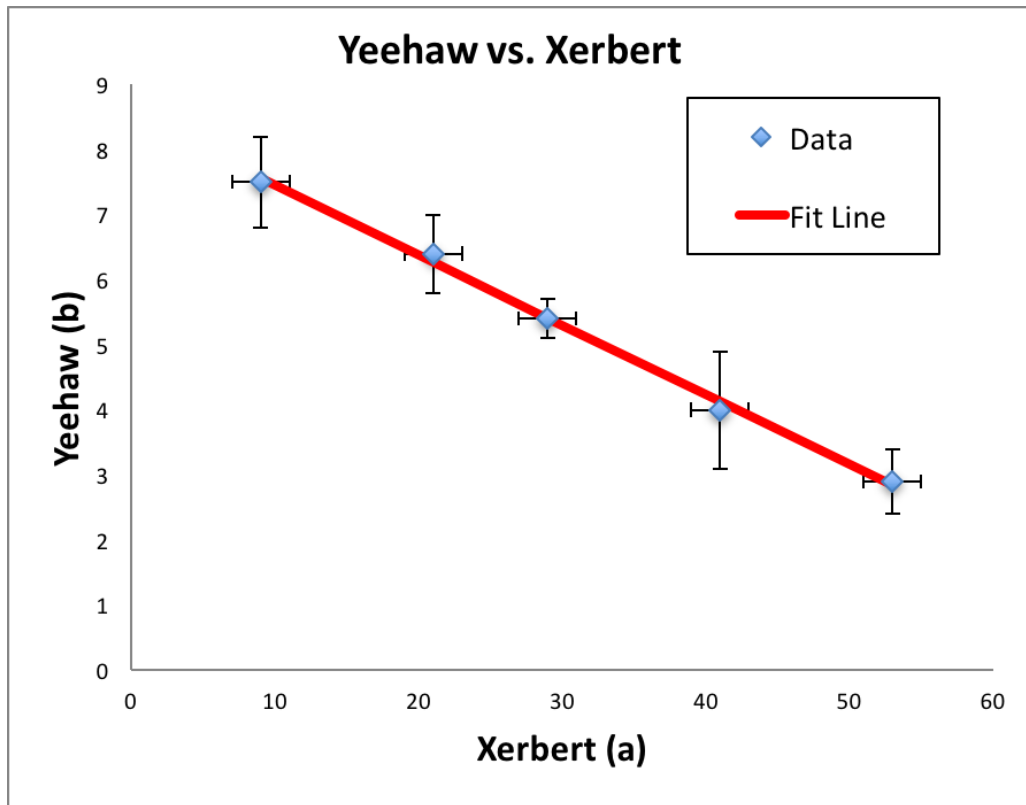
Figure 2.6: After running Solver, the original model is updated with the parameters from the successful minimization. The Solver "fit line" is a more successful model of the observed data.

# 3 PROCEDURE

The following exercise has been adapted from Problem 2.18 in **John R. Taylor's**, *An Introduction to Error Analysis (2nd Edition)*.

## 3.1 THE PHYSICS: PROJECTILE MOTION

Newton's Laws, and the subsequent equations of motion, can be used to describe the behavior of a vertically shot projectile. In particular, the theory can be used to define a mathematical relationship between the initial vertical velocity, $v$, and the maximum height achieved by the projectile, $h$. The relationship between these two measurable quantities is given by

$$v^2 = 2gh, \tag{3.1}$$

where $g$ is the acceleration due to gravity.

A group of students design an experiment to test the relationship given in Equation 3.1. A vertically aimed cannon is used to launch a projectile with varying initial velocities. The maximum height achieved by the projectile and the initial vertical velocity for each shot are measured and recorded to produce the data shown in Table 3.1.

| Trial | Height ($m$) | Velocity ($m/s$) |
|:-:|:-:|:-:|
| 1 | $0.4 \pm 0.1$ | $2.6 \pm 0.6$ |
| 2 | $0.8 \pm 0.1$ | $4.1 \pm 0.4$ |
| 3 | $1.4 \pm 0.1$ | $5.0 \pm 0.3$ |
| 4 | $2.0 \pm 0.1$ | $6.2 \pm 0.3$ |
| 5 | $2.6 \pm 0.1$ | $6.7 \pm 0.4$ |
| 6 | $3.4 \pm 0.1$ | $7.9 \pm 0.3$ |
| 7 | $3.8 \pm 0.1$ | $8.6 \pm 0.1$ |

Table 3.1: Experimental data.

## 3.2 ANALYSIS

This tutorial will walk you through a short analysis using the data presented in Table 3.1. The analysis techniques that are required here will be used routinely throughout the rest of the course.

1. Record the data listed in Table 3.1 in a well-formatted spreadsheet (see Section A.1).

2. Use a built-in spreadsheet function to compute the average uncertainty of the velocity measurements (see Section A.4). Be sure to label your result and report it using the appropriate number of significant figures (see Section A.3).

3. In a new column, compute the square of each velocity measurement, $v^2$. Remember: once an equation has been entered in a single cell, it can be easily copied and pasted to other cells.

   Be sure to clearly describe each column of your spreadsheet with a column title and the appropriate measurement units.

4. In a new column, compute the uncertainty associated with each squared velocity, $\delta(v^2)$. This process is called "error propagation"[*] (see Section A.3).

   The uncertainty $\delta(v^2)$ is found by taking the derivative of $v^2$ with respect to $v$ and multiplying by the uncertainty, $\delta v$:

$$\delta(v^2) = \frac{\partial(v^2)}{\partial v} \cdot \delta v \qquad (3.2)$$

$$\delta(v^2) = 2v \cdot \delta v. \qquad (3.3)$$

   Be careful to follow the rules for significant figures when reporting your answers.

5. Equation 3.1 is of the form:

$$y = ax, \qquad (3.4)$$

   where $y = v^2$ and $x = h$.

   Plot the squared initial velocity $v^2$ vs. the maximum height $h$. Be sure to format your plot as described in Section A.5.

6. Add error bars to the data points of your plot (see Section A.6).

7. Use `LINEST()` to determine the slope of the line and the associated uncertainty (see Sections A.7 and A.8). Be sure to report your result in standard form.

8. Construct your own linear model to describe the data as described in Section 2.2.

9. Use Solver to fit your linear model to the data. What physical quantity can be obtained from the slope of the trendline?

10. Compare the line fits obtained in steps 7 and 9 above. Do the two methods agree? Discuss any differences.

11. Quantitatively compare your result to the accepted value of $2g$. Assume that $g$ is known to be exactly $g = 9.81\,m/s^2$. Is the difference significant? Justify your answer.

12. Create a comparison plot like the ones shown in Taylor Chapter 2.4 to illustrate your result.

13. Discuss the possible (hypothetical) systematic uncertainty in the measured values of $v$ and $h$. In this case, you had nothing to do with the data collection and measurements, but based on your analysis of the data, it is still possible to look for systematic uncertainty.

   For example, what *systematic effect* in $v$ and $h$ would cause the measured result to differ from the accepted value? By "systematic effect" we mean: what if the velocity $v$ was measured to be larger than it actually was? What effect would this have on the calculated value of $2g$? What effect would you observe if it was measured to be smaller than it actually was? What about systematic effects in the measured height $h$?

---

[*]Error propagation quantifies the impact of measurement uncertainty on subsequent calculations – like the impact of $v \pm \delta v$ on calculating $v^2 \pm \delta v^2$ described above. Error propagation is one part of "Error Analysis". One of the learning goals of this course is to introduce you to the basic theory and mathematical tools used in error analysis.

14. **(Extra Credit):** The $\chi^2$ expression shown in Equation 2.4 is a simplified version compared to that given by Taylor (see Chapter 8.2, Equation 8.5). The simple version of the least-squares fit that is described above does not take into account the uncertainty on any given point. That is to say, all points are given the same weight in the fit. Suppose that we define the weight as:

$$w = \frac{1}{\sigma_{yi}^2}. \tag{3.5}$$

Develop a refined line fit that uses the weight function to perform a "weighted least-squares fit". Show that your refined fit gives data with smaller $y$-axis errors more influence on the fit result.

## 4  LAB REPORT

Before you leave the lab, submit your work to your T.A., via Blackboard. Normally, you will submit two files: a lab report (MS Word document) and an analysis spreadsheet (MS Excel spreadsheet). For this activity, a spreadsheet file (.xlsx) with your work and answers to the in-line questions will likely be sufficient. Consult your T.A. to determine how they would like you to submit your work.

# A  APPENDIX: SPREADSHEET REFERENCE

## A.1  RECORDING DATA IN A SPREADSHEET

Spreadsheets are designed to perform calculations on large collections of data. The first step is to record the data in a spreadsheet. In general, spreadsheets handle data best when they are recorded in columns (as opposed to rows). Columns must *always* be labeled with descriptive titles and the appropriate units. Measurement uncertainties are most conveniently recorded using a separate column. An example of some data recorded in a well-formatted spreadsheet is shown in Figure A.1.



| Trial (unitless) | Xerbert (a) | Xerb. Error (a) | Yeehaw (b) | Yeehaw Error (b) |
|---|---|---|---|---|
| 1 | 9 | 2 | 7.5 | 0.7 |
| 2 | 21 | 2 | 6.4 | 0.6 |
| 3 | 29 | 2 | 5.4 | 0.3 |
| 4 | 41 | 2 | 4 | 0.9 |
| 5 | 53 | 2 | 2.9 | 0.5 |

Figure A.1: An example of what your data might look like when recorded in a spreadsheet. Note that each column is clearly labeled with measurement units in parentheses.

## A.2  REFERENCING DATA

Data that are entered on a spreadsheet can be used and referenced in various ways. A specific cell can be referenced using the unique column-row identifier. For example B4 refers to the value recorded in the second column (B) and fourth row (4).

Groups of cells, in the same column or row, can be referenced as a *range*. For example, E2:E6 references all of the values listed in column E from row 2 through 8. Cells may also be referenced across columns and rows in an *array*. For example, D2:E6 specifies all of the numbers in the square bounded by D2, in the upper left and E6 in the bottom right.

## A.3  FORMATTING FOR THE APPROPRIATE SIGNIFICANT FIGURES

Measurements should always be reported with the appropriate number of significant figures. Two general rules should be followed when reporting a measurement:

1. Uncertainties should almost always be rounded to one significant figure.

2. The least significant figure should usually be the same order of magnitude (in the same decimal position) as the uncertainty.

Consider the data recorded in columns `D` and `E` in Figure A.1. The uncertainties in column `E` are all expressed using one significant figure according to Rule 1. Rule 2 states that each measurement should be expressed such that the least significant figure matches that of the uncertainty. This is true for all of the measurements except `D5`. The spreadsheet program has chosen to display `4.0` as 4 and now the least significant digits do not match the corresponding uncertainty in `E5` – violating Rule 2 above.

This behavior can be fixed by formatting the cells using the Format 〉 Cells menu (shown in Figure A.2) and setting the number of displayed decimal places to the appropriate value (in this case, 1).

For more information on significant figures, see Taylor's *"An Introduction to Error Analysis"*, Chapter 2.2
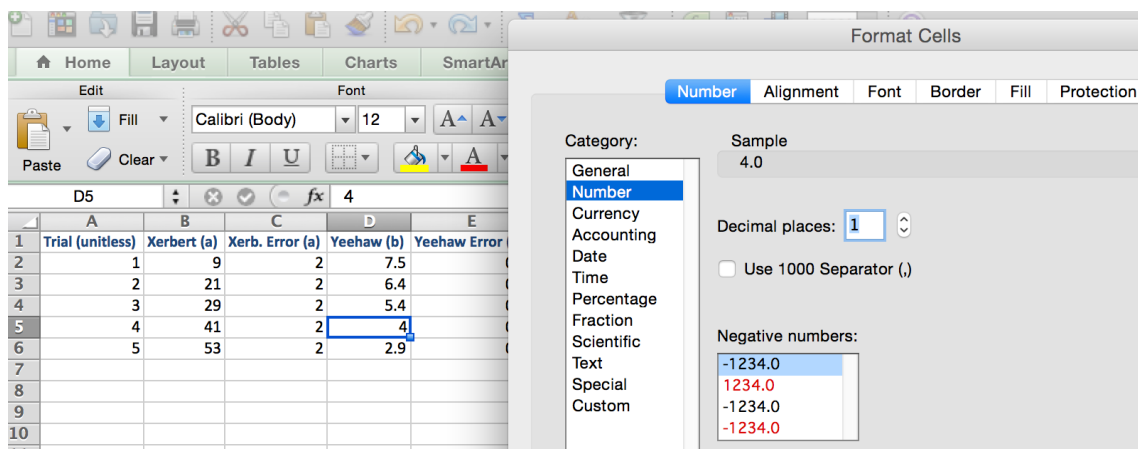


Figure A.2: Use the Format 〉 Cells menu to configure how cell contents are displayed. This is especially useful for formatting numbers to respect significant figure rules.

## A.4 CALCULATIONS AND FUNCTIONS

Cells that begin with "=" are interpreted by the spreadsheet as a calculation. Calculations can be anything from simple arithmetic to complicated expressions with multiple functions.

Arithmetic can be done by creating a custom equation. Standard arithmetic operators are interpreted as you would expect *e.g.* `+,-,*,/`. Exponents can be expressed by using the `^` operator.

Most spreadsheet programs include a large library of built-in functions that can be used to automate complicated calculations. A list of functions that a given spreadsheet program offers can be found by using the dedicated function button on the toolbar or by selecting the Insert 〉 Function menu option (Figure A.3).

Once an equation has been applied to a single cell, it can be copied and pasted to subsequent cells. Equations can also be expanded to multiple cells by using the "fill handle" on the lower-right corner of the highlighted cell.

Depending on the function being used, the input may be several comma-separated values and/or ranges. Functions may either output a single result or return an array of values. Built-in spreadsheet help is usually the best resource for learning how to correctly use a function. An example of the "Formula Builder" dialogue box is shown in Figure A.3

Functions that are commonly used in this course are listed in Table A.1. Other functions can be found using the function library included with your spreadsheet software. It is important to understand what the software is doing to compute the returned result.
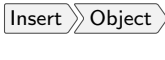


Figure A.3: Depending on your software, this is what the function wizard may look like. The function wizard can be accessed by clicking the button pointed to by the red arrow. Generally in the formula help window you will find a short description of the function and perhaps a link to more information on how the function calculates its advertised result.

| Function | Description |
|---|---|
| SUM(*range*) | Returns the sum of all arguments |
| AVERAGE(*range*) | Returns the average of all arguments |
| STDEV(*range*) | Calculates the standard deviation of a sample |
| POWER(*a*,*b*) | Returns $a^b$; $a$ raised to the power of $b$. |
| SQRT(*a*) | Returns $\sqrt{a}$. |
| LINEST(*y-range*,*x-range*,1,1) | Array formula that returns linear fit results. |

Table A.1: A short list of commonly used functions.

## A.5   VISUALIZING THE DATA – PLOTTING

Plots or graphs are often referred to as *charts* in spreadsheet software. To create a chart, you can click the "Chart Wizard" button on the tool bar or choose to insert a chart via the `Insert ⟩ Chart…` or `Insert ⟩ Object ⟩ Chart` menus.

The most common chart type used in this course is the "Marked Scatter" plot. An example of a poorly formatted scatter plot is given in Figure 2.1.

Plots that are submitted for this lab should meet the following formatting requirements:

- Graphs should be completely self explanatory – supporting text should not be required to explain what is plotted.

- The graph must have a clear title.

- All axes must be clearly labeled with units which should be spaced off and enclosed in parentheses. e.g. "xerbert (a)". Labels should follow the conventions, variables, and symbols used in your writeup.

- Include error bars, on both the X and Y axes when appropriate (Section A.6).

- When appropriate, include a trendline fit to the data points. The trendline equation and goodness-of-fit numbers should be printed on the plot (Section A.7).

- Include a legend *only* if multiple data sets are plotted on the same graph.

- In general, it is a good idea to remove the axis grid lines from the background of all plots.
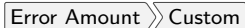
An example of a plot that meets all of the formatting requirements is shown in Figure A.4.

A checklist that can be used throughout this course to produce perfect plots is provided on the course Blackboard site under:
🗀 `Lab Assignments ▸ Reference Material ▸ Figure Formatting Guide`.

## A.6   ERROR BARS

Error bars can be added as a property of the chart or the "Data series" that is plotted. One way to add error bars is as follows:

- Add error bars by clicking on the chart then `Chart ⟩ Chart Layout ⟩ Error Bars ⟩ Error Bars Options`.

- In the "Format Error Bars" dialogue, specify the uncertainty range using `Error Amount ⟩ Custom`.

- Toggle between "X Error Bars" and "Y Error Bars" using the "Chart Elements" selection box in the "Current Selection" menu in the upper left-hand corner.
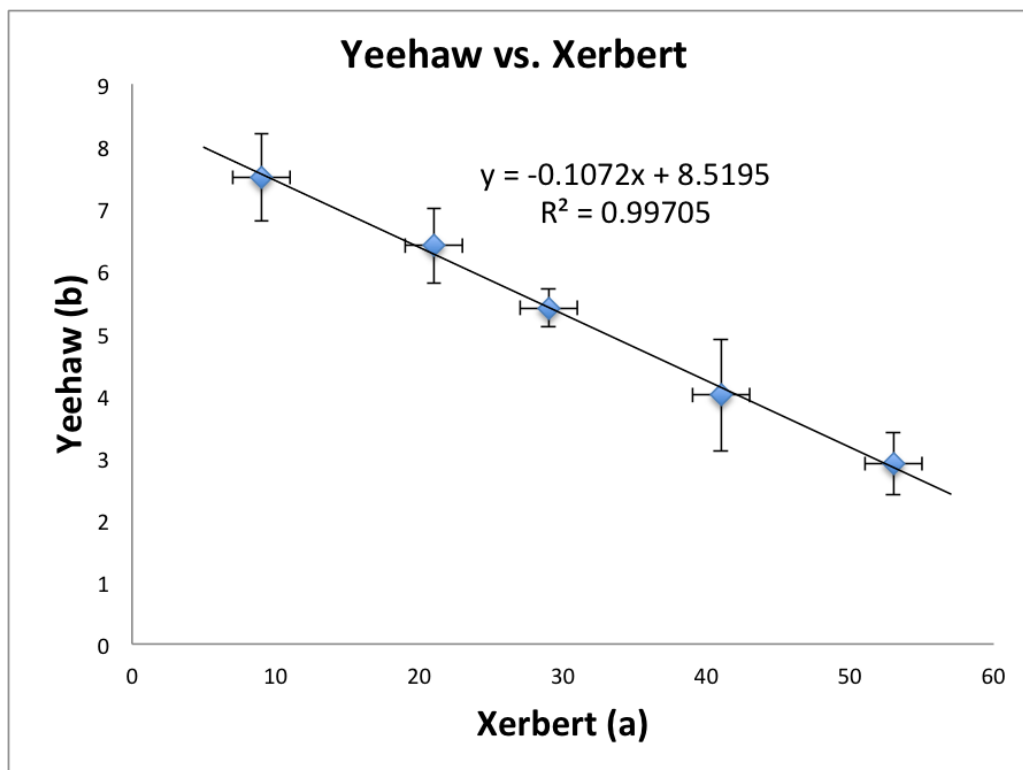
Figure A.4: An example of a well-formatted plot. The data are displayed with both x and y error bars and fit with a linear trendline. The equation for the line is shown on the chart along with the goodness-of-fit parameter, $R^2$.

## A.7 TRENDLINES

By eye, the data shown in Figure 2.1 appear to support a linear relationship. Line fits to data can determine the slope and tell us about the quality and nature of the relationship. Most spreadsheets refer to these "fits" as *trendlines*[†].

Trendlines are added to plots using a method that is similar to error bars. After clicking on the chart, Chart ⟩ Chart Layout ⟩ Trendline ⟩ Trendline Options... can be used to open the "Format Trendline" dialogue box.

Excel offers several different function shapes that can be used as trendlines. A few of the available options are listed in Table A.2. When a trendline is added to a chart, the fit results and fit quality should always be added to the plot. This can be done under the Options tab of the "Format Trendline" dialogue box. Select the "Display equation on chart" and "Display R-squared value on chart". R-squared is a parameter that describes the quality of the fit. For example, an $R^2$ of 1.0 would be a trendline where all of the data lie precisely on that line.

Figure A.4 shows a properly formatted trendline.

| Name | Functional Form |
|------|-----------------|
| Linear | $y = ax + b$ |
| Logarithmic | $y = a\ln(x) + b$ |
| Power | $y = ax^b$ |
| Exponential | $y = ae^{xb}$ |
| Polynomial (2nd Order) | $y = c_2 x^2 + c_1 x + c_o$ |

Table A.2: A short summary of some of the trendline functions that are available for fitting data in Excel.

## A.8 UNCERTAINTY ASSOCIATED WITH A TRENDLINE: LINEST()

With any fit to data, there is also an uncertainty associated with the fit parameters. Excel provides a function called LINEST() that returns information describing the parameters (slope and intercept), uncertainty, and quality of *linear* fits. The function LINEST() takes four arguments: a range of y-values, a range of x-values, and two options (that will always be "1" for our purposes). The output of LINEST() is a 2x3 cell array. Array functions must be entered in Excel following these steps:

1. In an empty cell, type the LINEST formula:

$$\texttt{=LINEST(<y-value range>, <x-value range>, 1,1)}.$$

2. Press return. A single value will appear in the cell where you typed the formula.

3. Highlight an array of cells that is 2 columns wide by 3 columns high – beginning with the cell where you typed the calculation.

4. Click in the formula toolbar (or use F2 (Windows) or ctrl + U (Mac)).

---

[†]It should be noted that the fitting done by most spreadsheet software is not very advanced and treats each point equally - ignoring the uncertainty associated with each data point.

5. Execute the array formula by pressing ⎡Ctrl⎤+⎡⇧⎤+⎡↵⎤ (Windows) or ⎡⌘⎤+⎡⇧⎤+⎡↵⎤ (Mac). Note that ⎡⇧⎤ denotes the "shift" key *not* the "up arrow" key.

You will now see that the 2x3 array of cells, that was highlighted in step 3 above, will now be populated with the results of the LINEST() calculation. The values returned by the LINEST() array are summarized in Table A.3.

The parameters returned by LINEST() are calculated using a "Vertical Least-Squares Fit". As such, LINEST() does not consider the error bars on the individual points when estimating the uncertainty for the fit parameters. A fit that also considers the uncertainty associated with each point is called a "Vertical Weighted Least-Squares Fit". See Chapter 8 of Taylor and Problem 8.9 specifically.

| Slope of the trendline | Intercept of the trendline |
|---|---|
| Estimated uncertainty for the slope | Estimated uncertainty for the intercept |
| $R^2$ goodness of fit parameter | – Not Relevant – |

Table A.3: Summary of the 2 column by 3 row array of values returned by LINEST().