

# *Introduction to Data Analysis and Plotting with Jupyter*

*Dustin Wheeler*

*Tuesday 26<sup>th</sup> January, 2021*

The goal of this exercise is to familiarize the student with basic data manipulation and plotting using Jupyter notebooks. Many of the concepts are general and can be translated to other software applications. By the end of this exercise, you should be able to

- perform basic manipulations on data columns (formatting, transformations)
- perform simple mathematical operations over a set of data
- plot columns of data in a clear and visually pleasing manner
- perform simple statistical operations on a set of data
- get a line-of-best-fit for multiple datasets
- import data from plain text files (CSV, TSV, DAT, etc.), like those created by many spectrometers and logging instruments

## *Introduction*

Why do we plot data? Edward Tufte may have said it best:

At their best, graphics are instruments for reasoning about quantitative information. Often the most effective way to describe, explore, and summarize a set of numbers—even a very large set—is to look at pictures of those numbers.

## *A Little Bit About Plots*

Plots (also called graphs or charts, though these are not *quite* the same thing) come in a variety of styles. They can be used to explore data as you are performing analyses or to convey results to a reader as part of a report. There are a variety of “standard” types which should cover your needs in this course. The main types we will use are:

*Scatter plot* Used to show correlation between two variables, each point generally represents an **independent** measurement.

Points are not **connected** by lines, lines are instead used to represent a mathematical model fitting the data.

*Line graph* Indicates continuity in the x dimension, often suggests a single collection session for all data on the “line” (e.g., a UV-Vis or NMR spectrum). Individual points are generally not shown. Lines should

never be smoothed (or “splined”), as this treatment no longer accurately represents the data.

Really, a line graph is just a special case of a scatter plot where the independent variable is continuous.

*Histogram* Represents the distribution of data in a series of “bins”, shows the frequency of a measured event. Closely related to density plots.

Here are the important parts of a plot, along with a brief description and guidelines for formatting.

*Title* Should be descriptive and represent *all* data in the plot.

*Axes* Should be clearly labeled (category, units, tick numbers). Should have an appropriate range (from slightly less than your minimum value to slightly more than your maximum value).

*Fit line* Should cover the whole spread of the data, but should not predict beyond the data range (you don’t know anything about measurements in that region).

*Legend* Should clearly show which set of data is which, trend lines should clearly correlate to a given data set.

### *Basic Plotting and Linear Fits*

Your first task is to take the data in the file `data.txt` and import it into a new Jupyter notebook. As the first step, please use this link to clone the files into your Jupyter directory: [https://jupyter.hcpchemlab.org/hub/user-redirect/git-pull?repo=https%3A%2F%2Fgithub.com%2Fmskblackbelt%2F357\\_plot\\_intro&urlpath=lab%2Ftree%2F357\\_plot\\_intro%2Fplot\\_intro\\_template.ipynb&branch=main](https://jupyter.hcpchemlab.org/hub/user-redirect/git-pull?repo=https%3A%2F%2Fgithub.com%2Fmskblackbelt%2F357_plot_intro&urlpath=lab%2Ftree%2F357_plot_intro%2Fplot_intro_template.ipynb&branch=main)<sup>1</sup> This link will direct your browser to the class JupyterHub and download a folder containing a number of resources you’ll need for this exercise.

The folder `raw_data` contains a few datasets in various formats. The file extension doesn’t always tell you what type of information is in a file and can sometimes lead you astray. For this exercise, the file `data.txt` contains a set of “tab-delimited” data, meaning that the values in a row are separated by tab characters. Another common format is “comma-delimited” or “comma-separated” data, such as the dataset in `anscombe.csv`. You will find that different instruments have different defaults for exporting data, but these two are the most common. Sometimes, a `.dat` extension might be used for a plain-text file to indicate ‘data’ without indicating the underlying format or delimiter type.

<sup>1</sup> Please do not copy and paste the text in this PDF into your input files, instead click on the link which should open your web browser. Invisible formatting characters are often included when copying from the PDF. Your web browser will misinterpret these characters and you will not be able to download the materials.

Open a new Jupyter template by clicking on the Template icon, then select 'pchem\_templates', and 'PChem\_Lab.ipynb' in the second box. This will open a new Jupyter notebook with a series of cells pre-populated with information to get you started. Give this notebook a title (as instructed in the first cell), then click on File in the JupyterLab menubar, select "Rename Notebook...", and give the notebook a title following the instructions in the first cell. In the notebook, execute the code cell under the Library Import section, then create a new cell under the Data import section. Import the data from `data.txt` by using the import function of the 'pandas' library:

```
1 data = pd.read_csv("raw_data/data.txt", sep="\t")
```

This will take the contents of `data.txt` and turn it into an object called a "DataFrame" with the name `data`. The `sep="\t"` argument is required because our data is "tab-separated", rather than "comma-separated", which is the default (as suggested by the function name `read_csv`). When in doubt, specify the field separator when you import data with python and the pandas library.

To view your new object, just type `data` and press `Shift + Enter` to evaluate the cell. By default Jupyter notebooks will display the results of the final line of a cell. You should end up with two columns of data with 21 data points in each column. Once you have the data columns, create a new scatter plot from the data. The pandas library automatically incorporates basic plotting functions from the matplotlib library, so data in a DataFrame can be plotted by using the `plot()` function on a given DataFrame. In your case, you need only type `data.plot()` and evaluate the cell to see a plot of the data.<sup>2</sup> While you are exploring data, the default appearance is fine, but you'll probably want to start customizing the appearance of plots when you get them ready for publication.

Notice that the plot isn't plotting column X against column Y, but instead is plotting X and Y against the *index* (the bold column of numbers on the left side of the displayed data). There are two ways to correct your plot. The first is to import the data and define the index column during the import, like so:

```
1 data = pd.read_csv("raw_data/data.txt", sep="\t",
2   index_col="X")
```

This method is useful if you know the structure of the data ahead of time (and when you have a clear index column, such as a time index, wavelength, or frequency). The second way to fix the plot is to define the X and Y data in the arguments of the `plot()` function. The `plot()` function

For more information, type `help(pd.DataFrame)` to get more information on the DataFrame object. Alternately, in the Jupyter notebook, you can press `Shift + Tab` at the end of a command to get help on the object (i.e., when your cursor is at the end of `pd.DataFrame` or your new object `data`.)

They only display the results of intermediate lines of a cell if the `print()` command is used.

<sup>2</sup> *n.b.*: pandas is designed first as a tool for finance and business operations. The default settings for plots are often inappropriate for scientific plots so some tweaking may be necessary.

takes a number of arguments (you can see them quickly by pressing `Shift` + `Tab` inside of the parentheses), including arguments for features like plot type (`kind="text"`), plot title (`title="text"`), and axes labels (`xlabel="text"`, `ylabel="text"`).

With this in mind, create a scatter plot (`kind="scatter"`) of Y and a function of X, give the plot a title (“X vs. Y” is fine) and label each axis (“X”, “Y”).

### *Polishing Plot Appearance*

When creating more complex plots, once the basics are complete, remove any unnecessary “chart junk” to prepare your plot for publication. This includes things like background shading (`matplotlib` is pretty good about this), grid lines (only necessary if you’re trying to manually read data off the chart, we don’t need this with the ubiquity of computers), and the border surrounding the full plot area (should only need axes to indicate the origin or provide an anchor for tick marks). If there is a legend (useful when there are multiple datasets present), it should usually be a part of the plot area, not pushed off to the side. Don’t waste valuable data area on something like a legend than can be nestled in one of the empty parts of your plot. Finally, while color seems like a nice feature, plots shouldn’t rely on color to convey meaning. There’s always a chance that readers won’t have access to color printouts, or that the reader may be colorblind. To keep plots accessible, use symbols, line weights, and dashing to differentiate your datasets. The lesson to learn from this paragraph is that simplicity is key to conveying your message. The sole purpose of your figures is to help the reader understand the message. Anything that distracts from that is harmful to your cause.

Now that the plot has been cleaned up, we’d like to add a trendline. The `seaborn` library is a visual analysis package built on top of `matplotlib` with an emphasis on data exploration and friendly figures. To plot our data with `seaborn`, do the input the following:

```
1 sns.lmplot(data=data, x="X", y="Y", ci=None)
```

This creates a scatterplot using the X and Y columns from `data`, fits it with a linear model (“lm”), and disables the confidence interval (CI) in the plot. A linear fit should be sufficient for this data. The `seaborn` library isn’t meant for complex customization or complex

### *Fitting Data*

Now that you have a trendline, it’s time to get some information about it. The `numpy` library is a fundamental library in the numeric python

community. It contains a vast array of functions and structures that are used by *many* other libraries. One useful function in `numpy` is the `polyfit()` function, which calculates a polynomial of best fit using a linear least-squares method. Examples can be seen on the linked page, but we will run through the simplest example, a straight line.

To use the `polyfit()` function, you need to specify the data (“X” and “Y”) and the degree of the polynomial (1). All other arguments are optional, and just giving these three arguments will return the coefficients of the polynomial from the highest power to the lowest ( $c_n, c_{n-1}, \dots, c_1, c_0$ ). The optional argument `cov=True` will return the *covariance matrix*, the diagonal values of which are the variances for each coefficient.

These (or, more precisely, their square roots) help you to see how “good” the returned fit is. Please read the help page to learn what statistics are returned and in what form.

### Anscombe’s Quartet

Repeat the process of importing data with the file `anscombe.csv`. Name the new dataframe `anscombe` and display it by calling it alone in the last line of your cell.

Now we want to calculate and print a linear fit and corresponding error for each of the four sets of data in `anscombe`. The columns in `anscombe` should be named “x1”, “y1”, “x2”, “y2”, etc. For convenience, I’ve created a small loop for you to use on your data. This allows us to perform repetitive operations without worrying about being consistent in our typing from operation to operation (if it’s correct for the first set of data, it will be correct for each subsequent one).

If you’re interested in doing more statistical modeling in Python, the `statsmodel` package offers a number of descriptive statistics and flexible models.

There are two options here: you will need to change the delimiter to a comma (`delim=","`) for this file, or, since the default delimiter for `pd.read_csv()` is a comma, you can omit the `delim=","` argument altogether (it is implicitly defined).

Just like we did with the dataframe `data` in the previous section.

```
1  ### We will make a loop to make this less
   ↳ repetitive (and reduce our chances of making an
   ↳ error)
2  fits = []
   for x, y in
   ↳ (('x1', 'y1'), ('x2', 'y2'), ('x3', 'y3'), ('x4', 'y4')):
4      fit_coeff, fit_cov = np.polyfit(anscombe[x],
   ↳ anscombe[y], 1, cov=True)
      fit_err = np.sqrt(np.diag(fit_cov))
6      fits.append([fit_coeff, fit_err])

8  # This operation makes my point regarding this data
   ↳ more clear, as you'll see
   with np.printoptions(precision=2):
10     print("Fit, Std. Err.")
      for x in fits:
12         print(x)
```

Clearly, the data sets are different... You can see that from the numbers in the table. So what's going on? To learn more, plot each of the datasets in a separate plot. Discuss your observations with your lab partner, and figure out the best way to deal with information like this.

The key takeaway from this exercise is that you should always plot your data, even if it's only a quick glance on the side of your spreadsheet! This data was taken from [an article written by F.J. Anscombe in 1973](#) arguing for the necessity of graphs in statistical analysis. It is such a valuable dataset that the `seaborn` library is installed with a copy of the data (`sns.load_dataset("anscombe")`).

## Column Manipulation and Plotting Spectra

Start a new section titled “Plotting Spectra” and import the file `ftir_data.csv`. Print out the top of the dataframe using the command `df.head()`, where `df` is the name of your dataframe. Notice that we have data in the form “Wavenumber ( $\text{cm}^{-1}$ )”, “Transmittance (%)”.

There is a very simple syntax in `pandas` for making a new column:

```
1 df["new col name"] = func(df["old col name"])
```

In this case, `func()` just represents the transform that we are applying to the original column (multiplication, addition, inversion, etc.). Columns in dataframes will broadcast functions across every item in the column by default, so making a line-by-line transformation (e.g., doubling a column) is as easy as `df["B"] = 2 * df["A"]`. Using this syntax, create new calculated columns in our dataframe so we can plot the same spectrum as “Wavelength (nm)” vs “Absorbance (A)”.<sup>3</sup>

To get you started converting wavenumber to wavelength, make the following transformation on the “Wavenumber” column:

```
1 df["Wavelength"] = 1/(df["Wavenumber"]*10**(-2)) * 10**6
```

You should be able to derive this equation quickly on your own...<sup>4</sup>

Perform the same steps to take the column for “Transmittance” and transform it into Absorbance values. Remember that Absorbance is  $\log_{10}(1/T)$  while Transmittance is  $I/I_0$ .

Once you have created the new columns, create two plots. The first should plot wavenumber against transmittance, while the second will plot wavelength against absorbance. Because these are spectra (data collected in a single experiment), it is appropriate to plot them as line plots, *without* plotting the individual points. Make sure that your plot reflects this.

Make sure your plots are labeled in all the appropriate locations and

<sup>3</sup> What are the relevant units for %T and A?

Notice that Python requires the use of “`***`” for exponentiation. This is because the caret character (`^`) is often used for other commands. This convention is common in programming languages.

<sup>4</sup> Hint: label the units and work through the dimensional analysis.

that your spectra look clean and presentable.

Congratulations, you have the basics of plotting and calculation using Jupyter! You are expected to take this information and build upon it as you write reports through the semester. All of these techniques can be translated to other applications, some of which might be more appropriate for the type of research being performed in a research group. You should be able to transfer your basic skills to any other software with ease, and the tips on presentability are universally applicable.

### *Alternative Software*

While Microsoft Excel may be the most common piece of software for plotting, that doesn't mean it's the best. You may find that you desire more control over your plots, or a simpler interface, or the ability to manipulate your data in other manners. If so, you might want to look at the following applications. Some are free, some are available through a CUNY site license, and some are paid (and expensive). This list is not comprehensive, but it does give you an idea of the variety available to you should you be interested.

*LibreOffice Calc* (Free) Part of the free and open source LibreOffice suite, designed as a full featured alternative to Microsoft Office.

<https://www.libreoffice.org/discover/calc/>

*Apple Numbers* (Free) Apple's entry into the spreadsheet category. Works nearly as well as Excel, cleaner interface, but macOS/iOS-only and can't always do some of the more advanced statistics.

<https://www.apple.com/numbers/>

*Google Sheets* (Free) Like all Google products, this is web-only (no offline data analysis here!). Designed to be a clone of Excel.

<https://docs.google.com/spreadsheets/>

*GraphPad Prism* (CUNY site license available) Statistical analysis software for Windows and macOS. Capable of repeated analysis on huge datasets of life/social science data. Falls short in spectral analysis tools.

<https://www.graphpad.com>

*Origin by OriginLab* (no site license available) Professional data analysis and plotting software for scientists and engineers. Heavy focus on physical sciences and engineering. Only available on Windows.

<https://www.originlab.com>

*Igor Pro by WaveMetrics* (no site license available) Very similar to Origin, but available for macOS and Windows.

<https://www.wavemetrics.com/products/igorpro>

*Mathematica* (site license), *Matlab* (site license), *JupyterLab* (free) These are really programming environments, though they are designed for use by scientists. Mathematica and Matlab are commercial products, while Jupyter is an open-source option that uses the **JULia**, **PYThon**, and **R** languages. They all have entry-output interfaces so that live interaction can be done with data. While the learning curve is steep, the payout for learning a programming language will last a lifetime. We will be using the Python language later in the semester for some computational chemistry work.

<http://www.wolfram.com/mathematica/>

<https://www.mathworks.com/products/matlab.html>

<https://jupyter.org>