

COMP 364: Computer Tools for Life Sciences

Using libraries: NumPy & Data visualization with Matplotlib

Christopher J.F. Cameron and Carlos G. Oliver

Key course information

Midterm

- ▶ how was it? too hard? too easy?
- ▶ if you could not write the midterm
 - ▶ please contact Chris or Carlos ASAP

Quiz #5

- ▶ will be available October 30th, 2017
- ▶ access closes at 11:59:59 pm on the same day
- ▶ MC questions will cover topics from the last two weeks

Assignment #3

- ▶ coming soon!

Whitespace cheatsheet

Common whitespace characters:

- ▶ ' ' or ' '- space
- ▶ '\t' - a tab
- ▶ '\n' - newline character ('\r\n' in Window environments)
- ▶ '\r' - carriage return (move to the beginning of the line)
- ▶ '\f' - form feed (advance to next page)
- ▶ '\s' - any whitespace

Data

Let's think back to last Friday...

- ▶ randomly chose 10 genes and two time points
- ▶ learned how to label genes based on a similarity measure
- ▶ in our case, gene expression across two chosen time points
- ▶ then by cluster analysis, we grouped genes together

```
1 start,end = [170 ,240]
2 gene_names = ["YMR274C", "YOR263C", "YLR371W",
3               "YDL120W", "YKL096W", "YHR192W",
4               "YMR046C", "YDR038C", "YMR172W",
5               "YDL045C"]
6 obs = [(-0.24, -0.03), (-0.22, 0.25),
7         (0.06, 0.1), (-0.02, -0.57),
8         (0.66, -0.73), (0.24, 0.34),
9         (0.17, -0.75), (0.59, 0.27),
10        (-0.08, 0.18), (0.07, 0.23)]
11 cluster_labels = [0, 0, 1, 2, 2, 1, 2, 1, 0, 1]
```

Cluster analysis - *k-means*

In our last lecture, we ran the *k-means* algorithm

- ▶ using an implementation from Python's SciPy module
 - ▶ `scipy.cluster.vq.kmeans2()`

k-means algorithm:

Step 1 - randomly choose k centroids

Step 2 - for each data point, assign it to the nearest centroid

Step 3 - recalculate each cluster's new centroid

Step 4 - repeat step 2

Step 5 - if data points have changed clusters, go to step 3
if no data points have been reassigned, stop

Example - <https://www.youtube.com/watch?v=5I3Ei69I40s>

Matplotlib

Having numerical data is useful, but...

- ▶ sometimes it is more practical to visualize data in a plot

To do this we are going to use Python's **Matplotlib** module

- ▶ allows for 2D and 3D plotting
- ▶ provides useful functionality for Python to work like MATLAB
 - ▶ another useful programming language

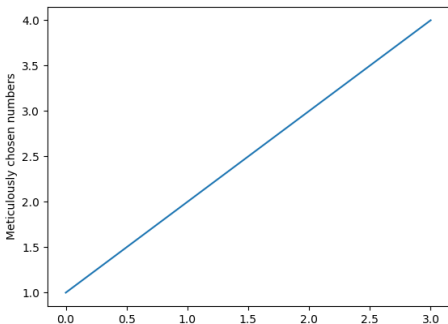
Let's start by importing the module:

```
1 import matplotlib.pyplot as plt
```

Matplotlib #2

Let's create a simple plot:

```
1 import matplotlib.pyplot as plt
2 plt.plot([1,2,3,4])
3 plt.ylabel("Meticulously chosen numbers")
4 plt.show() # displays figure
```



Matplotlib #3

Why does the x-axis ranges from 0-3 and the y-axis from 1-4?

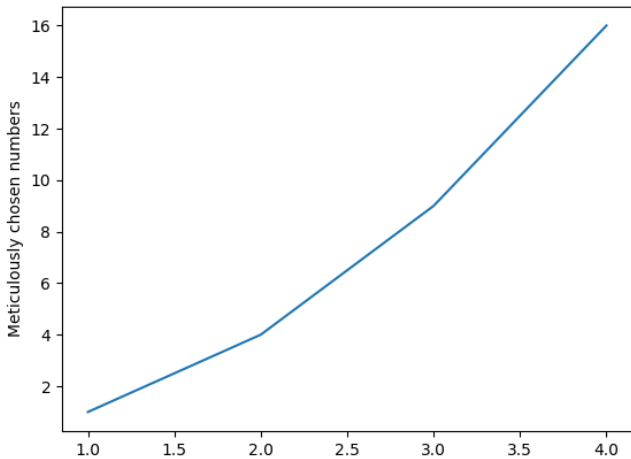
If you provide a single list or array to the **.plot()** command

- ▶ matplotlib assumes it is a sequence of y values
- ▶ and automatically generates the x values for you
- ▶ python ranges start with 0
- ▶ the default x vector has the same length as y but starts with 0

To plot x versus y:

```
1 plt.plot([1, 2, 3, 4], [1, 4, 9, 16])
```

Matplotlib #4



Jupyter notebooks

To use matplotlib in the Jupyter notebooks

- ▶ add **'%matplotlib inline'**
- ▶ before the import statement of the module

```
1 %matplotlib inline
2
3 import matplotlib
4 import matplotlib.pyplot as plt
```

You could also change the IPython kernels

- ▶ but this is outside the scope of the course
- ▶ <https://stackoverflow.com/questions/19410042/how-to-make-ipython-notebook-matplotlib-plot-inline>

Useful matplotlib functions

To save a figure to file use **.savefig()**

- ▶ https://matplotlib.org/devdocs/api/_as_gen/matplotlib.pyplot.savefig.html

When producing multiple plots

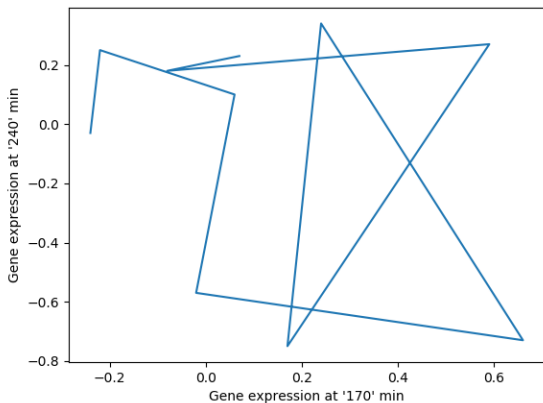
- ▶ make sure to close the previous one using **.close()**
- ▶ https://matplotlib.org/devdocs/api/_as_gen/matplotlib.pyplot.close.html

Let's try plotting our gene expression data from before

Plotting gene expression data

```
1 # split obs by x and y vals
2 x_vals,y_vals = zip(*obs)
3 plt.plot(x_vals,y_vals)
4 # set x and y labels
5 plt.xlabel("Gene expression at '"+str(start)+"' min")
6 plt.ylabel("Gene expression at '"+str(end)+"' min")
7 # write figure to file
8 plt.savefig("./gene_figure.png")
9 # close matplotlib figure
10 plt.close()
```

Plotting gene expression data #2



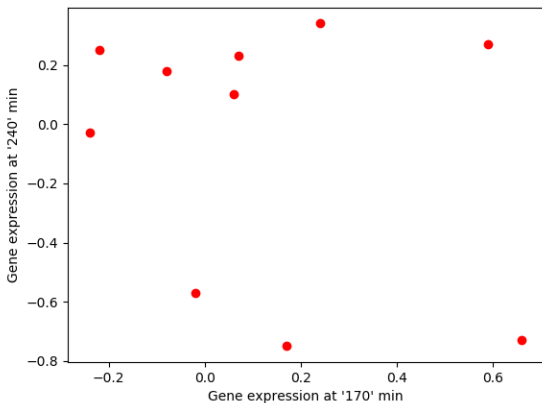
That doesn't quite look right

- ▶ let's look at the matplotlib API

Plotting gene expression data #3

```
1 # split obs by x and y vals
2 x_vals,y_vals = zip(*obs)
3 plt.plot(x_vals,y_vals,"ro")
4 # set x and y labels
5 plt.xlabel("Gene expression at '"+str(start)+"' min")
6 plt.ylabel("Gene expression at '"+str(end)+"' min")
7 # write figure to file
8 plt.savefig("./gene_figure_2.png")
9 # close matplotlib figure
10 plt.close()
```

Plotting gene expression data #4



That's better, but how can we colour data points by cluster?

- ▶ let's look at the matplotlib API

Colours in matplotlib

Matplotlib functions can handle many different colour codes:

1. RGB
 - ▶ '(0, 0, 0)' - black
2. hex RGM
 - ▶ '#0F0F0F' - black
3. character
 - ▶ 'b' - blue
 - ▶ 'k' - black
 - ▶ 'r' - red
4. many more...

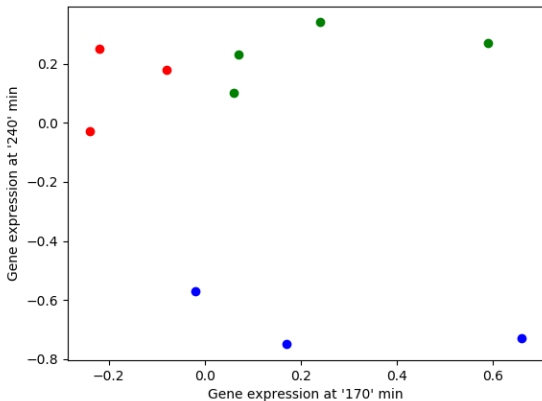
For more information, see:

https://matplotlib.org/api/colors_api.html

Plotting gene expression data #5

```
1 # define set of colors to use
2 colors = ["r", "g", "b"]
3 # get set of labels used
4 labels = set(cluster_labels)
5 # iterate over labels and plot subsets
6 for label in labels:
7     color = colors[label]
8     vals = []
9     for i,cluster in enumerate(cluster_labels):
10        if(cluster == label):
11            # index obs for cluster label
12            vals.append(obs[i])
13    x_vals,y_vals = zip(*vals)
14    # plot current cluster, based on label
15    plt.plot(x_vals,y_vals,"o",color=color)
```

Plotting gene expression data #6



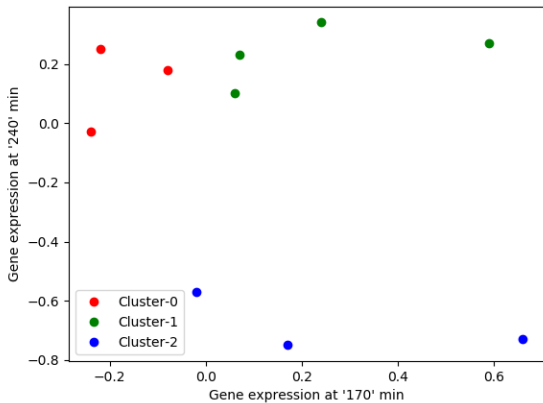
Even better, let's add a legend

- ▶ let's look at the matplotlib API

Plotting gene expression data #7

```
1  # iterate over labels and plot subsets
2  for label in labels:
3      color = colors[label]
4      vals = []
5      for i,cluster in enumerate(cluster_labels):
6          if(cluster == label):
7              # index obs for cluster label
8              vals.append(obs[i])
9      x_vals,y_vals = zip(*vals)
10     # plot current cluster, based on label
11     plt.plot(x_vals,y_vals,"o",
12              color=color,
13              label="Cluster-"+str(label))
14 plt.legend(loc="best")
```

Plotting gene expression data #8



Gene expression across a time course

Let's try something a little more difficult

- ▶ begin by reading in all expression data for chosen genes

```
1 data_dict = {}
2 with open("./Spellman.csv", "r") as f:
3     header = f.readline().rstrip().split(",")
4     time_points = [int(val) for val in header[1:]]
5     for line in f:
6         line = line.rstrip().split(",")
7         gene_name = line[0]
8         # if gene was randomly chosen
9         if(gene_name in gene_names):
10             exps = [float(val) for val in line[1:]]
11             data_dict[gene_name] = exps
```

Gene expression across a time course #2

Okay, now we have a Python dictionary containing:

- ▶ keys that represent gene names
- ▶ values containing a gene's expression across the time course

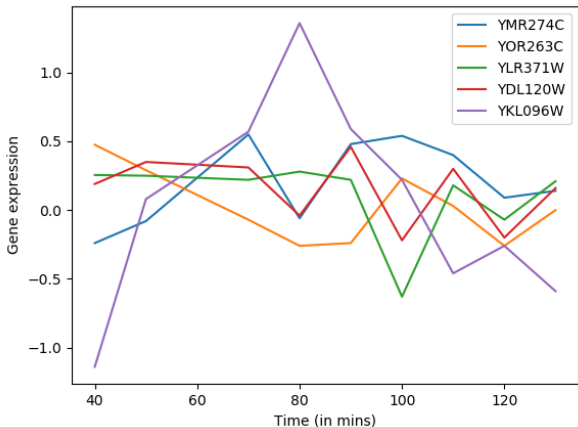
```
1 gene_name = "YMR046C"
2 print(data_dict[gene_name])
3 #prints:
4 #[0.0, -0.09, 0.145, 0.38, 0.1, 0.02, 0.03, 0.43,
5 # -0.56, 0.38, -0.32, 0.23, -0.5, 0.17, 0.12, 0.21,
6 # -0.09, 0.09, -0.17, 0.11, -0.75, -0.11, 0.045]
```

How can we display the expression of these genes as a line plot?

Gene expression across a time course #3

```
1 # iterate over genes
2 for gene_name in gene_names[:5]:
3     # plot corresponding time points and
4     # gene expression
5     plt.plot(time_points[:10], data_dict[gene_name][:10],
6              label=gene_name)
7 plt.ylabel("Gene expression")
8 plt.xlabel("Time (in mins)")
9 plt.legend(loc="best")
10 plt.savefig("./gene_figure.png")
11 plt.close()
```

Gene expression across a time course #4



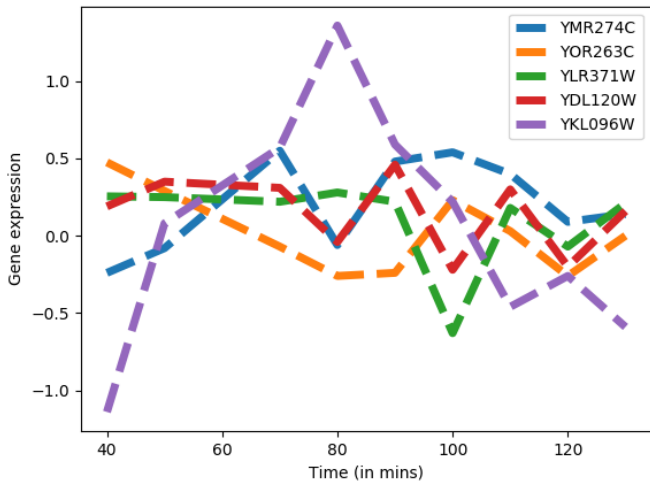
Can I make the lines thicker? Different style?

- ▶ let's look at the matplotlib API

Gene expression across a time course #5

```
1 # iterate over genes
2 for gene_name in gene_names[:5]:
3     # plot corresponding time points and
4     # gene expression
5     plt.plot(time_points[:10], data_dict[gene_name] [:10],
6              label=gene_name, ls="--", lw=5)
7 plt.ylabel("Gene expression")
8 plt.xlabel("Time (in mins)")
9 plt.legend(loc="best")
10 plt.savefig("./gene_figure.png")
11 plt.close()
```

Gene expression across a time course #6



Next time in COMP 364

Continuing our plunge into Python's Matplotlib module

Integrating Python's **NumPy** module

- ▶ a fundamental package for scientific computing, which contains:
 - ▶ a powerful N-dimensional array object
 - ▶ useful linear algebra and random number capabilities

NumPy API:

<https://docs.scipy.org/doc/numpy-1.13.0/reference/>