# IS262B Daniels Final Project WIP

June 4, 2020

This Jupyter Notebook (JN) is meant to serve as a template for one way in which researchers can publish their data alongside their research findings. In addition to a JN like this one, the method of data publication for which we are advocating also involves the use of MyBinder. Because the UCLA Data Science Center is working to develop UCLA's instance of Dataverse, this workflow is intended for research data that is going to be ingested into UCLA's Dataverse instance Section ??. As such, the three major components needed for this workflow are Jupyter Notebook (preferably accessed with a package manager like Anaconda), MyBinder, and Dataverse.

By utilizing JN, we are able to break down code into manageable chunks, and we are also able to easily annotate, explain, or format our code using Markdown. The end goal is for documents like this to be ingested into Dataverse alongside their respective dataset(s). This will add transparency and interoperability to publications that require, encourage, or permit the publication of research data in addition to the research findings themselves. It can also help to further explicate the methodology employed to analyze or work with the dataset(s), giving authors a way to more easily explain their thought process and workflow. Additionally, because Jupyter Notebook can be used by scholars to do actual research and not just as a way to publish data after the fact, it's possible for researchers to use this workflow without adding a significant amount to their workload.

Which brings us to the final note about this document. The Gapminder data that we are working with is not part of a specific research project. This is appropriate in this instance because we are developing a workflow that involves using a brand new feature in Dataverse: integrating MyBinder and Jupyter Notebook so that users can run code for a given dataset directly from Dataverse. As such, it would not make sense to use "live" research data as fodder for experimentation, since reusing old data allows us the time and flexibility to thoroughly test our workflow for a variety of potential use cases.

In a "real" use case, our JN would contain the code and accompanying markdown text that produced the results or findings of the researcher(s). In our case, because we lack research results or some kind of "product", we are instead doing some basic data cleaning and data visualization with the data. Additionally, over the course of developing this workflow we have identified two different approaches for writing this kind of JN. The first method, which we employed, is to provide explanations for an audience that has little to no experience with Python. The second approach is much more likely to be employed, which is to write the JN with just basic explanations, assuming your audience has at least intermediate experience with Python or another programming language Jupyter supports.

Although the approach that we used is admittedly more work, in theory this work does not necessarily need to be done by the PI, for example. Instead, explicating code for a given dataset could be done by experts like the Data Science Center staff in the UCLA Library in exchange for partial authorship, for example. Regardless, there can exist a clear division of labor between the "actual" work that researchers are most interested in doing and this kind of work. As academic research

increasingly revolves around one or more digital datasets, as well as the new digital research tools necessary to work with datasets, the need to provide supplemental documentation like this JN becomes more and more apparent.

## 1 Gapminder Data Visualization

#### 1.0.1 Overview

This Jupyter Notebook (JN) uses Python to visualize existing GDP data from the Gapminder Foundation. These visualizations rely on two Python modules, matplotlib and pandas. All code necessary to run these visualizations is contained below, including code that loads the two requisite Python modules for you. To begin, we will load in the appropriate Python modules, matplotlib and pandas. The syntax below means that we can "call" or reference both matplotlib and pandas with the shorthand plt and pd, respectively. There is strong convention within the Python community to use plt and pd as the standard shorthand for their respective modules and we follow that convention here.

```
[90]: %matplotlib inline import matplotlib.pyplot as plt import pandas as pd
```

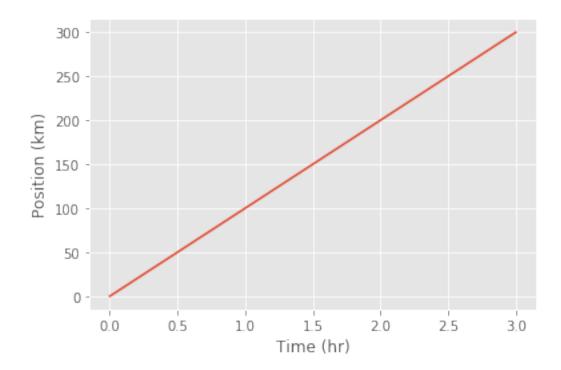
We will now plot a very simple graph with some arbitrary data that we create. The arbitrary data below consists of two lists, "time" and "position," each of which contains a list of four numbers. We plot these numbers in a simple graph to confirm that our matplotlib module is loaded and functioning.

Note the syntax plt.plot, plt.xlabel, plt.ylabel, etc. This code is referencing the matplotlib module using our designated shorthand, plt.

```
[91]: time = [0, 1, 2, 3]
    position = [0, 100, 200, 300]

    plt.plot(time, position)
    plt.xlabel('Time (hr)')
    plt.ylabel('Position (km)')
```

[91]: Text(0, 0.5, 'Position (km)')



Now we can start working with the Gapminder data. Here we have included the code for loading the pandas module, even though we've already loaded it, to better illustrate some of the functionality of the pandas module. Pandas allows us to load a dataset into our Python program. In this case, we are loading up gapminder\_gdp\_oceania.csv and are assigning this dataset to the variable "data". This variable will function in much the same way as the shorthand we used to access matplotlib and pandas. Note that although variable names can be anything, it is considered good practice to use variable names that relate to, or in some way otherwise elicit, the dataset that you are assigning to the variable. Choosing helpful variable names will not only allow others to more easily read and understand your code down the line, but will also help you stay organized as you write your code.

Our code was written to look for our dataset in the same directory as this JN, in a folder called "data". However, if your CSVs are in a different directory, you can simply input your correct directory address in place of data/gapminder\_gdp\_oceania.csv – just be sure to include quotation marks around your directory, for example:

```
pd.read\_csv('c:\users\yourname\downloads\data\gapminder\_gdp\_oceania.csv', index\_col='country')
```

You'll notice we are also assigning the column named 'country' as our index. An index refers to a position within an ordered list. In our case, the name of the country makes the most sense.

```
[92]: import pandas as pd

# switching to reading in tab delimited file (the format dataverse converts csv

→to)

#data = pd.read_csv('data/gapminder_gdp_oceania.csv', index_col='country')
```

country	gdpPercap_1952	gdpPercap_1957	gdpPercap_1962	gdpPercap_1967	\
country Australia New Zealand	10039.59564 10556.57566	10949.64959 12247.39532	12217.22686 13175.67800	14526.12465 14463.91893	
country	gdpPercap_1972	gdpPercap_1977	gdpPercap_1982	gdpPercap_1987	\
Australia New Zealand	16788.62948 16046.03728	18334.19751 16233.71770	19477.00928 17632.41040	21888.88903 19007.19129	
country	gdpPercap_1992	gdpPercap_1997	gdpPercap_2002	gdpPercap_2007	
Australia New Zealand	23424.76683 18363.32494	26997.93657 21050.41377	30687.75473 23189.80135	34435.36744 25185.00911	

With print(data) we can see that our CSV has been properly loaded and assigned to the "data" variable. Next we'll want to clean and standardize the data to make it easier for us to work with. The first thing we'll do is simplify column names. In addition to the 'country' column we have columns that refer to different dates. Fortunately these date columns all follow the same naming convention, which will make it that much easier to clean the data. Common to all date columns is that the column name always starts with "gdpPercap\_XXXX" where "XXXX" refers to the four digit year. For clarity, we're going to remove the first portion of each column name using the .strip() function, leaving behind only the year.

If you're curious, once you have your csv loaded into your 'data' variable, you can check to see what the columns are named by using the following code: print(data.columns) You should see a list of strings like this: Index(['gdpPercap\_1952', 'gdpPercap\_1957', 'gdpPercap\_1962', 'gdpPercap\_1962', 'gdpPercap\_1962', 'gdpPercap\_197', 'gdpPercap\_1982', 'gdpPercap\_1982', 'gdpPercap\_1987', 'gdpPercap\_1992', 'gdpPercap\_1997', 'gdpPercap\_2002', 'gdpPercap\_2007'], dtype='object')

Those names are unnecessarily long, and we are only interested in the dates. In programming parlance, we are going to remove each instance of a string, or set of characters. In our case, this means going through each column in the dataset and removing all instances of our string. Surprisingly, this is very simple to do and can be done with a single line of code:

```
[93]: years = data.columns.str.strip('gdpPercap_')
```

The syntax here is relatively simple. We are telling the program to create the variable years, to which we are going to access our other variable, data. However, from that variable we want only the column headers, which we communicate with the .columns addition. Furthermore, because the .strip() command only works on strings, we need to make sure that everything that is being parsed during our .strip() command is a string. We do that by simply adding a .str to our line of code. In plain English, data.columns.str.strip('gdpPercap\_') says from the variable data, access each column, turn the value into a string, and then strip all strings that are an exact match of

gdpPercap\_

Now we have a variable named years that contains a list of the corrected column names. We can check with:

```
[94]: print(years)
```

Now we need to rename each column with the corresponding date in our list. During this step we also need to convert these numbers from strings into integers. Converting to integers is necessary for our matplotlib module to work.

The following code is reassigning values to the columns in our dataset while simultaneously converting these values into integers:

```
[95]: data.columns = years.astype(int)
```

data.columns will iterate over each column name just like before, but this time instead of reading these values we are going to assign new values to each. Specifically we are iterating through our list, years, and assigning each of those values to our columns as integers. Now if we were to print the columns, instead of the clunky list we saw before we should see a list of dates, exactly like we saw in our years variable:

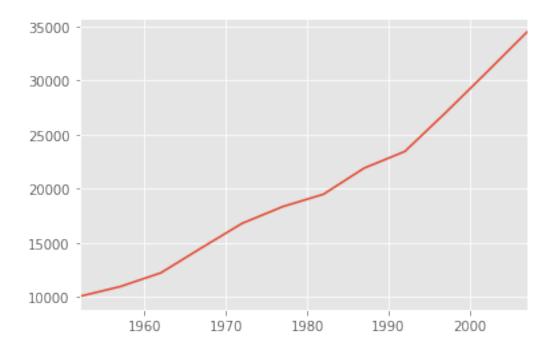
```
[96]: print(data.columns)
```

Now that our data's been sufficiently cleaned we can start plotting data directly from our variable, data, which contains what is known as a Pandas dataframe. Let's plot a quick line graph:

For more information on Pandas dataframes, follow this link: https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.html

```
[97]: data.loc['Australia'].plot()
```

[97]: <matplotlib.axes. subplots.AxesSubplot at 0x7fc610b02e10>



At a glance, the graph suggests that our data is cleaned sufficiently for us to start graphing. We can take a look ourselves by looking at the head or the tail of the dataset; that is to say, the first five rows or the last five rows, respectively. To do that, we use the .head() or .tail() function, like so:

The first five rows of our data are:

[98]:	country	1952	1957	1962	1967	1972	\
	Australia New Zealand	10039.59564 10556.57566	10949.64959 12247.39532	12217.22686 13175.67800	14526.12465 14463.91893	16788.62948 16046.03728	
	country	1977	1982	1987	1992	1997	\
	country	10004 10751	10477 00000	01000 00000	00404 74400	04007 00457	
	Australia	18334.19751	19477.00928	21888.88903	23424.76683	26997.93657	
	New Zealand	16233.71770	17632.41040	19007.19129	18363.32494	21050.41377	
		2002	2007				
	country						
	Australia	30687.75473	34435.36744				
	New Zealand	23189.80135	25185.00911				

You'll notice that in our case, we only have two rows in total, so if you were to run the .tail() function instead, you would get the same result. But what if you wanted to reverse the table by transposing

the rows to the columns and vice versa? In other words, instead of having gdpPercap\_<year> in the column across the top, we have Australia, New Zealand and the gdpPercap years as columns? We can simply use the T function

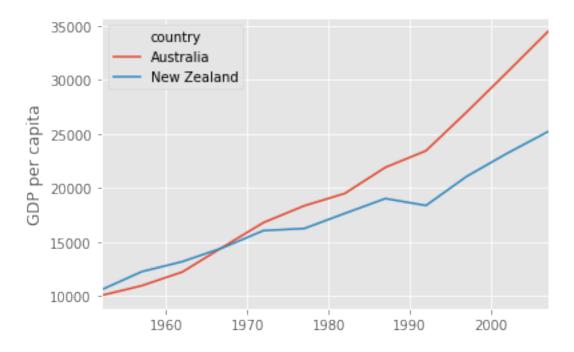
#### [99]: data.T

```
[99]: country
                 Australia
                            New Zealand
      1952
               10039.59564
                             10556.57566
      1957
               10949.64959
                            12247.39532
      1962
               12217.22686
                            13175.67800
      1967
               14526.12465
                            14463.91893
      1972
               16788.62948
                            16046.03728
      1977
               18334.19751
                            16233.71770
      1982
               19477.00928
                            17632.41040
      1987
                            19007.19129
               21888.88903
      1992
               23424.76683
                            18363.32494
      1997
               26997.93657
                            21050.41377
      2002
               30687.75473
                            23189.80135
      2007
               34435.36744
                            25185.00911
```

Now let's plot it.

```
[100]: data.T.plot()
   plt.ylabel('GDP per capita')
```

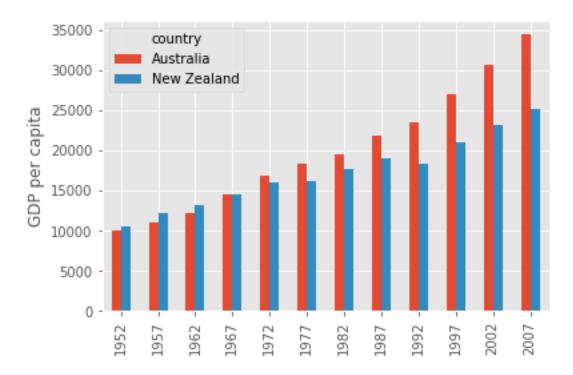
[100]: Text(0, 0.5, 'GDP per capita')



There are many preconfigured styles for graphs. Next we'll use one called "ggplot" to make a bar graph. Note the x-axis ticks are now vertical, allowing for more dates to fit into the same amount of space.

```
[101]: plt.style.use('ggplot')
   data.T.plot(kind='bar')
   plt.ylabel('GDP per capita')
```

```
[101]: Text(0, 0.5, 'GDP per capita')
```

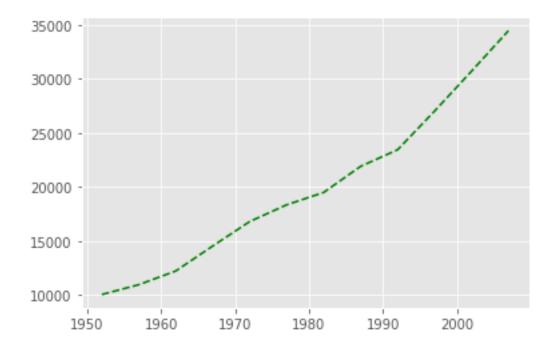


So far we've only been using Python's most basic graphing functions. For more detailed or otherwise more sophisticated graphs, we'll need to use the plot function from the matplotlib module.

Up until now we've relied on Python's default plotting behavior to produce our graphs. Modules like matplotlib allow us to exercise more granular control over our data. As an example, next we'll make two variables that we will use to fill in our x- and y- axis data: years will be our x-axis values and gdp\_australia will be our y-axis values. The three comma-separated values that are contained within the .plot parentheses are not only how we control which of our data we want to feed into our graph, but also allow us to customize the visualizations. In the following example, our .plot() function has input values for our graph's x-axis, y-axis, and the style of our line (in this case, a green dashed line represented by g-):

```
[102]: years = data.columns
gdp_australia = data.loc['Australia']
plt.plot(years, gdp_australia, 'g--')
```

[102]: [<matplotlib.lines.Line2D at 0x7fc660949d90>]



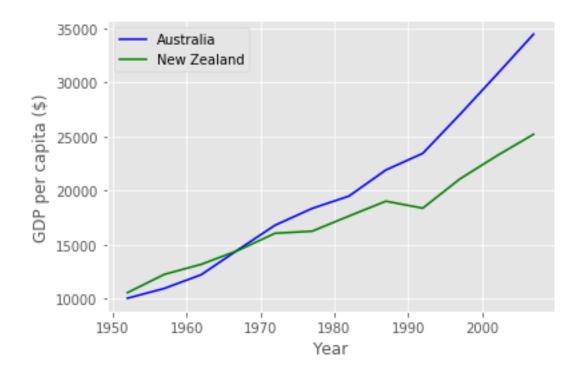
Matplotlib enables us to plot multiple sets of data together on the same graph. In addition to assigning our x-axis values to a variable, this time we'll also assign values to our y-axis. Let's also add some labels to each axis, as well as a legend to label our lines. Note that it appears as though we only gave instructions on the location of the legend. In reality, matplotlib is able to populate the legend because we gave each of our line plots a label in our .plot() function. Moreover, specifying the location of the legend is optional; matplotlib will try and place the legend in a suitable place on its own by default.

```
[103]: # Select two countries' worth of data.
gdp_australia = data.loc['Australia']
gdp_nz = data.loc['New Zealand']

# Plot with differently-colored markers, add x- and y-axis labels.
plt.plot(years, gdp_australia, 'b-', label='Australia')
plt.plot(years, gdp_nz, 'g-', label='New Zealand')

# Create legend.
plt.legend(loc='upper left')
plt.xlabel('Year')
plt.ylabel('GDP per capita ($)')
```

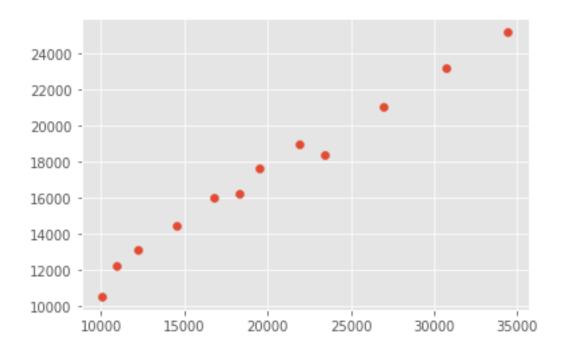
[103]: Text(0, 0.5, 'GDP per capita (\$)')



In addition to line graphs, matplotlib can produce many other types and styles of graph, such as a scatter plot. To make one all we have to do is use the .scatter() function in much the same way as we used the .plot() function. In fact, even though we are currently working in a different code block, the variables that we've created up until this point are still valid and accessible. Let's use the same variables as before to populate our scatter plot's x- and y-axis:

```
[104]: plt.scatter(gdp_australia, gdp_nz)
```

[104]: <matplotlib.collections.PathCollection at 0x7fc6211bd910>

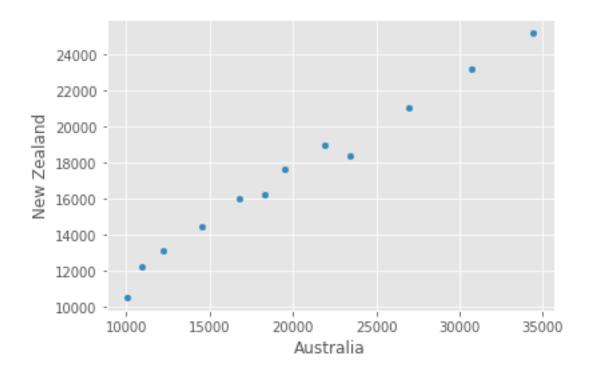


Programming languages like Python more often than not allow for multiple different solutions to the same problem. As a small example, we'll use a different method for labeling our axes.

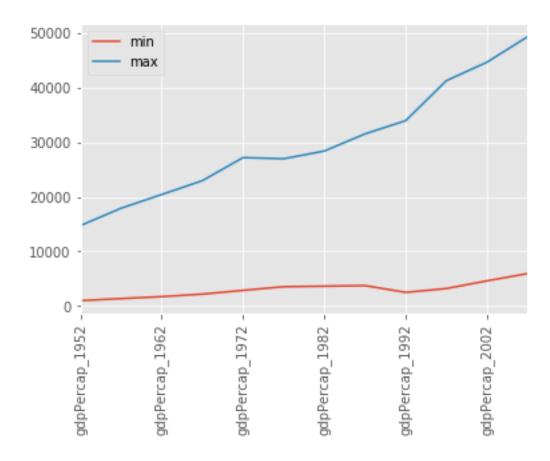
Notice that with this code block, we are writing all of our code on just a single line using method chaining.

```
[105]: data.T.plot.scatter(x = 'Australia', y = 'New Zealand')
```

[105]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fc610ccdbd0>

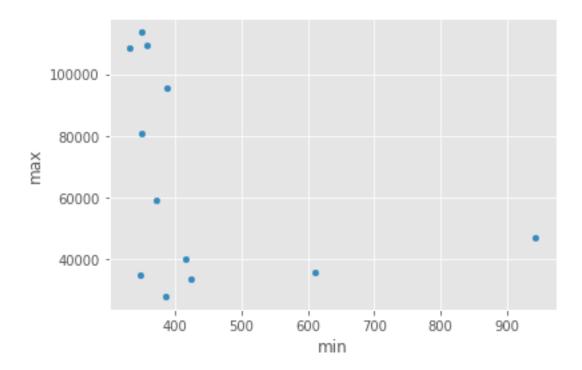


Many functions in Python can be chained together using what's called a method chain. Method chaining slims down your code, not only saving you space but also often making your code easier for other's read and understand. Let's import some new data and plot a graph using method chaining and the .min() and .max() functions.



Visualizing data like this can help researchers make quick, broad analyses of a given dataset, thereby helping to reveal interesting relationships in the data. Such analyses also often reveal new lines of inquiry to scholars. For example, what (if any) relationship do you see between the minimum and maximum GDP per capita data among Asian countries for each year in the dataset? We can produce a scatter plot to help us answer this question:

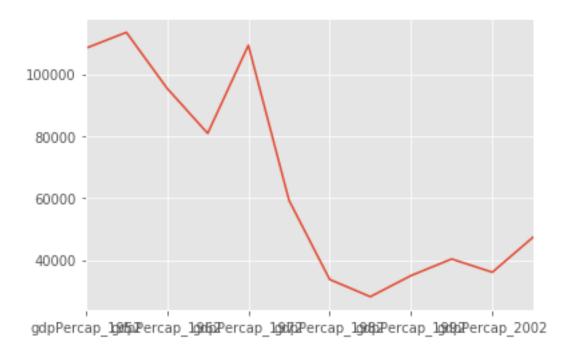
[107]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fc610ac7e10>



It appears as though there are no particular relationships, nor any obvious correlations between the min and max GDP values year by year. The fortunes of Asian countries, according to this data, do not rise and fall together. However, the graph does show us that there is a lot more variability among the maximum values than the minimum values. To investigate, let's take a look at those max values and their corresponding index. (Recall that we set our index to be the unique values in the 'country' column.)

https://pandas.pydata.org/pandas-docs/version/0.15.2/generated/pandas.DataFrame.plot.html

[108]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fc6607b0ed0>

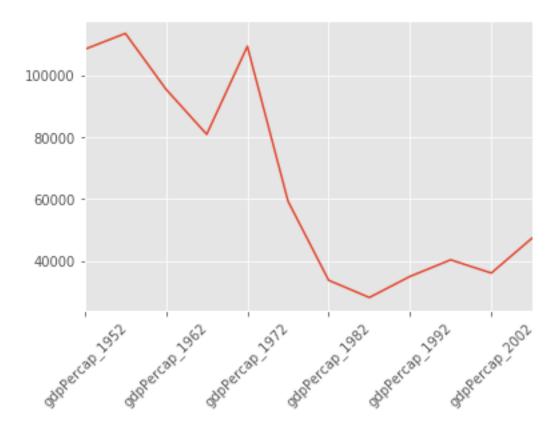


Right away you'll notice our labels (x-ticks) are illegible. The .plot(), like other functions in Python, accepts different parameters. That is to say, we can modify the output of our function, or alternatively control how the functions inputs or interprets information. For example, when we read our csv file we've been using the index\_col parameter to set our index values. We will use the rot parameter to rotate our x-ticks 45 degrees so that they are legible. As you may have guessed, the rot parameter takes an integer as its argument, or input. For a full list of parameters and what they do, see this documentation:

https://pandas.pydata.org/pandas-docs/version/0.15.2/generated/pandas.DataFrame.plot.html

```
[109]: data_asia.max().plot(rot = 45)
```

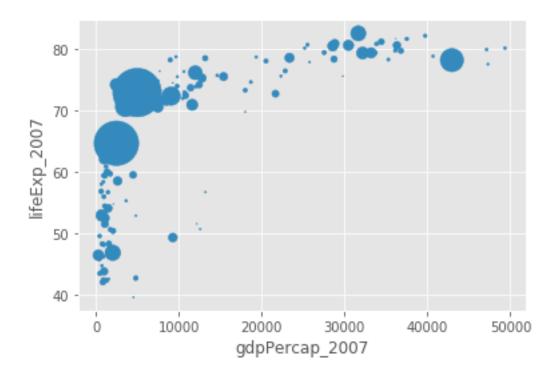
[109]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fc670de8bd0>



The graph reveals a precipitous drop after 1972, possibly alluding to some kind of global, geopolitical event.

Our Gapminder datasets also includes life expectancy data. The following short code block will show the correlation between a country's life expectancy and its GDP for 2007, normalizing marker size by population using the scatter plot's s argument:

[110]: <matplotlib.axes.\_subplots.AxesSubplot at 0x7fc670e54d50>



Lastly, let's download some of these graphs by using the .savefig() function:

```
[111]: plt.savefig('my_figure.png')
```

<Figure size 432x288 with 0 Axes>

In this example, your graph will be saved as a PNG file called "my\_figure" and will be saved to your current working directory. If you are unsure of your current working directory, you can simply type "pwd", which stands for Print Working Directory, and then run it and Python will print your current working directory. You can change the file format by changing the file extension to one of the other available formats: pdf, ps, eps, and svg.

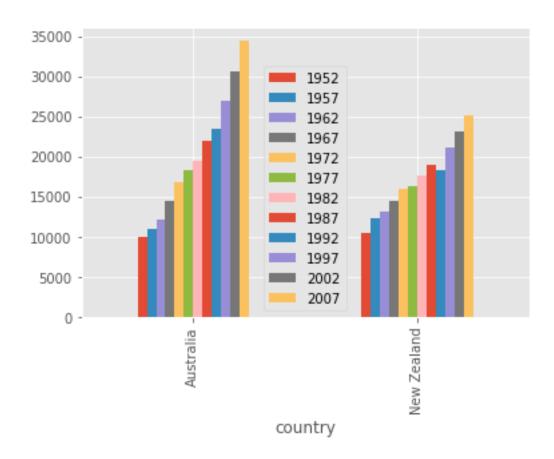
```
[112]: pwd
```

### [112]: '/Users/dougdaniels/Documents/GitHub/is262b\_final'

It is often the case that when you work with dataframes, your data is generated and plotted to screen in a single line (method chaining), and the plt.savefig() function doesn't seem viable. A possible workaround is to assign a reference to your figure to a local variable using plt.gcf(), and to then use the .savefig() function on your variable:

```
[113]: fig = plt.gcf() # get current figure
data.plot(kind='bar')
fig.savefig('my_figure.png')
```

<Figure size 432x288 with 0 Axes>



1. While working on this project we found a bug in how mybinder works with the UCLA Dataverse instance. We are working with the mybinder group to identify and fix the bug. For this project, we used the Harvard Dataverse instance.