



(<https://cognitiveclass.ai>).

Pie Charts, Box Plots, Scatter Plots, and Bubble Plots

In []:

#by Christopher Harrison

Introduction

In this lab session, we continue exploring the Matplotlib library. More specifically, we will learn how to create pie charts, box plots, scatter plots, and bubble charts.

Table of Contents

1. [Exploring Datasets with *pandas*](#)
2. [Downloading and Prepping Data](#)
3. [Visualizing Data using Matplotlib](#)
4. [Pie Charts](#)
5. [Box Plots](#)
6. [Scatter Plots](#)
7. [Bubble Plots](#)

</div>

Exploring Datasets with *pandas* and Matplotlib

Toolkits: The course heavily relies on *pandas* (<http://pandas.pydata.org/>) and *Numpy* (<http://www.numpy.org/>) for data wrangling, analysis, and visualization. The primary plotting library we will explore in the course is *Matplotlib* (<http://matplotlib.org/>).

Dataset: Immigration to Canada from 1980 to 2013 - [International migration flows to and from selected countries - The 2015 revision](http://www.un.org/en/development/desa/population/migration/data/empirical2/migrationflows.shtml) (<http://www.un.org/en/development/desa/population/migration/data/empirical2/migrationflows.shtml>) from United Nation's website.

The dataset contains annual data on the flows of international migrants as recorded by the countries of destination. The data presents both inflows and outflows according to the place of birth, citizenship or place of previous / next residence both for foreigners and nationals. In this lab, we will focus on the Canadian Immigration data.

Downloading and Prepping Data

Import primary modules.

In [1]:

Let's download and import our primary Canadian Immigration dataset using *pandas* `read_excel()` method. Normally, before we can do that, we would need to download a module which *pandas* requires to read in excel files. This module is **xlrd**. For your convenience, we have pre-installed this module, so you would not have to worry about that. Otherwise, you would need to run the following line of code to install the **xlrd** module:

```
!conda install -c anaconda xlrd --yes
```

Download the dataset and read it into a *pandas* dataframe.

In [2]:

Data downloaded and read into a dataframe!

Let's take a look at the first five items in our dataset.

In [3]:

Out[3]:

	Type	Coverage	OdName	AREA	AreaName	REG	RegName	DEV	DevName	198
0	Immigrants	Foreigners	Afghanistan	935	Asia	5501	Southern Asia	902	Developing regions	1
1	Immigrants	Foreigners	Albania	908	Europe	925	Southern Europe	901	Developed regions	
2	Immigrants	Foreigners	Algeria	903	Africa	912	Northern Africa	902	Developing regions	8
3	Immigrants	Foreigners	American Samoa	909	Oceania	957	Polynesia	902	Developing regions	
4	Immigrants	Foreigners	Andorra	908	Europe	925	Southern Europe	901	Developed regions	

5 rows × 43 columns

Let's find out how many entries there are in our dataset.

In [4]:

(195, 43)

Clean up data. We will make some modifications to the original dataset to make it easier to create our visualizations. Refer to *Introduction to Matplotlib and Line Plots* and *Area Plots, Histograms, and Bar Plots* for a detailed description of this preprocessing.

In [5]:

data dimensions: (195, 38)

Visualizing Data using Matplotlib

Import Matplotlib.

In [6]:

Matplotlib version: 3.1.1

Pie Charts

A `pie chart` is a circular graphic that displays numeric proportions by dividing a circle (or pie) into proportional slices. You are most likely already familiar with pie charts as it is widely used in business and media. We can create pie charts in Matplotlib by passing in the `kind=pie` keyword.

Let's use a pie chart to explore the proportion (percentage) of new immigrants grouped by continents for the entire time period from 1980 to 2013.

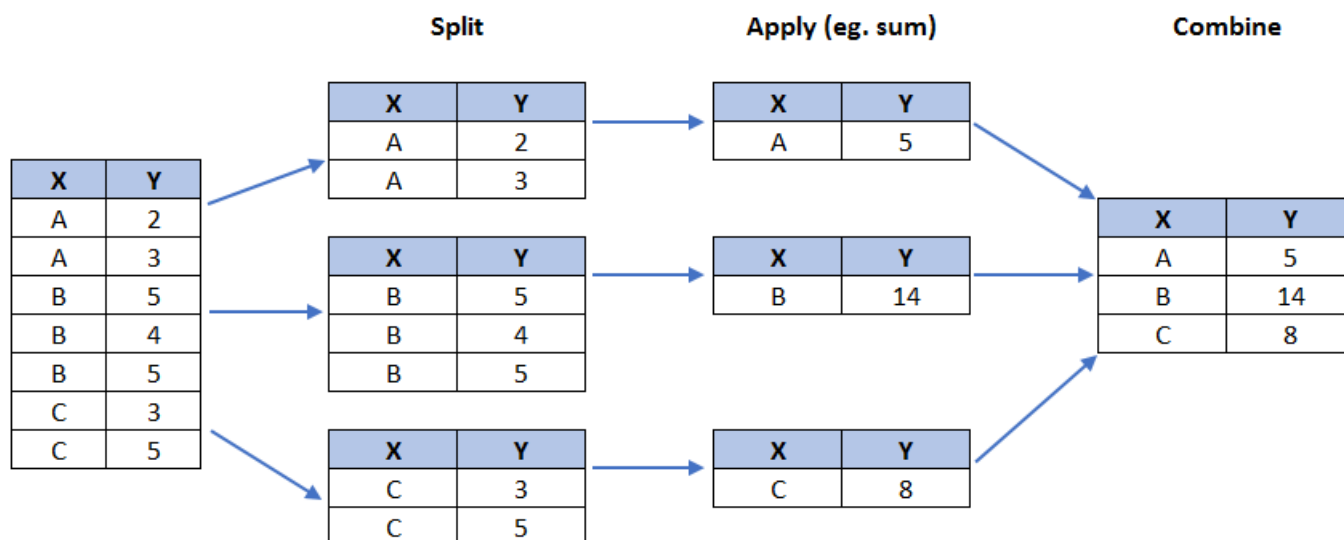
Step 1: Gather data.

We will use `pandas` `groupby` method to summarize the immigration data by `Continent`. The general process of `groupby` involves the following steps:

1. **Split:** Splitting the data into groups based on some criteria.
2. **Apply:** Applying a function to each group independently:

```
.sum()
.count()
.mean()
.std()
.aggregate()
.apply()
.etc..
```

3. **Combine:** Combining the results into a data structure.



In [7]:

```
<class 'pandas.core.groupby.generic.DataFrameGroupBy'>
```

Out[7]:

	1980	1981	1982	1983	1984	1985	1986	1987	1988	1989	...	200
Continent												
Africa	3951	4363	3819	2671	2639	2650	3782	7494	7552	9894	...	2752
Asia	31025	34314	30214	24696	27274	23850	28739	43203	47454	60256	...	15925
Europe	39760	44802	42720	24638	22287	20844	24370	46698	54726	60893	...	3595
Latin America and the Caribbean	13081	15215	16769	15427	13678	15171	21179	28471	21924	25060	...	2474
Northern America	9378	10030	9074	7100	6661	6543	7074	7705	6469	6790	...	839

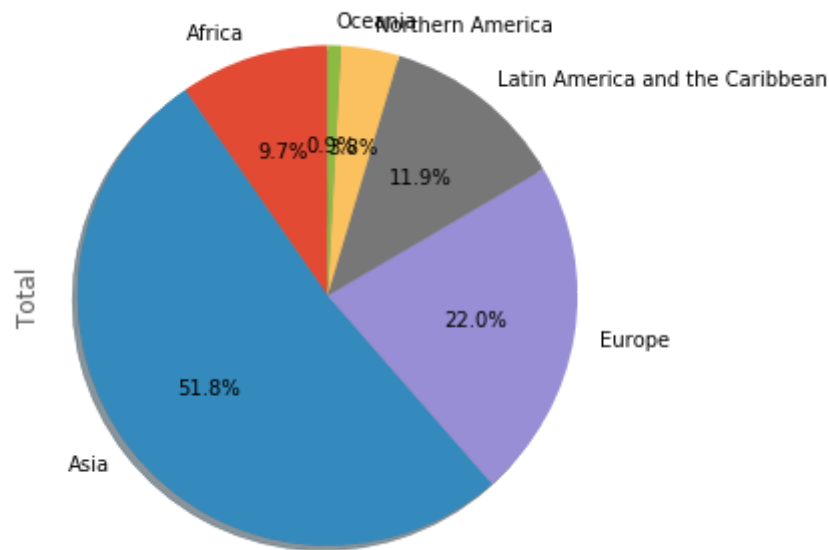
5 rows × 35 columns

Step 2: Plot the data. We will pass in `kind = 'pie'` keyword, along with the following additional parameters:

- `autopct` - is a string or function used to label the wedges with their numeric value. The label will be placed inside the wedge. If it is a format string, the label will be `fmt%pct`.
- `startangle` - rotates the start of the pie chart by angle degrees counterclockwise from the x-axis.
- `shadow` - Draws a shadow beneath the pie (to give a 3D feel).

In [8]:

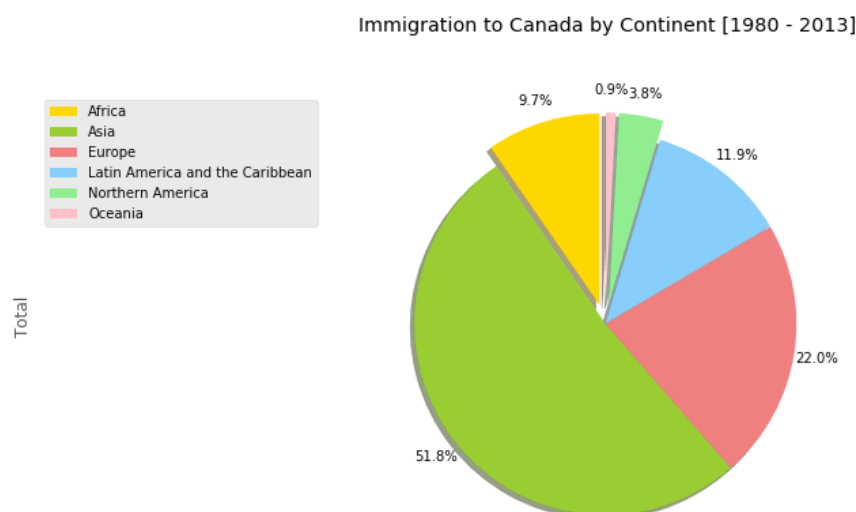
Immigration to Canada by Continent [1980 - 2013]



The above visual is not very clear, the numbers and text overlap in some instances. Let's make a few modifications to improve the visuals:

- Remove the text labels on the pie chart by passing in `legend` and add it as a separate legend using `plt.legend()`.
- Push out the percentages to sit just outside the pie chart by passing in `pctdistance` parameter.
- Pass in a custom set of colors for continents by passing in `colors` parameter.
- **Explode** the pie chart to emphasize the lowest three continents (Africa, North America, and Latin America and Caribbean) by passing in `explode` parameter.

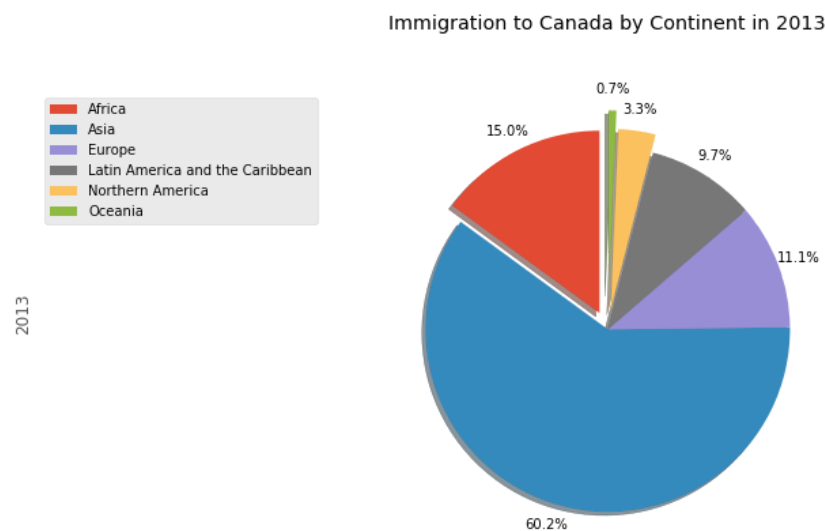
In [9]:



Question: Using a pie chart, explore the proportion (percentage) of new immigrants grouped by continents in the year 2013.

Note: You might need to play with the explore values in order to fix any overlapping slice values.

In [15]:

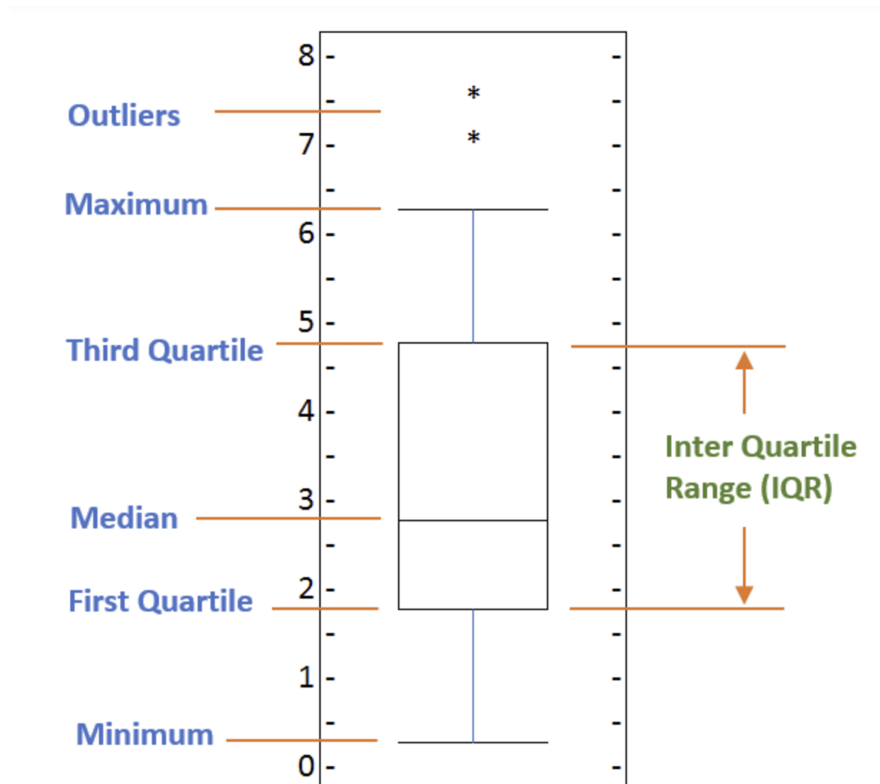


Double-click **here** for the solution.

Box Plots

A box plot is a way of statistically representing the *distribution* of the data through five main dimensions:

- **Minimum:** Smallest number in the dataset.
- **First quartile:** Middle number between the minimum and the median.
- **Second quartile (Median):** Middle number of the (sorted) dataset.
- **Third quartile:** Middle number between median and maximum.
- **Maximum:** Highest number in the dataset.



To make a box plot, we can use `kind=box` in `plot` method invoked on a *pandas* series or dataframe.

Let's plot the box plot for the Japanese immigrants between 1980 - 2013.

Step 1: Get the dataset. Even though we are extracting the data for just one country, we will obtain it as a dataframe. This will help us with calling the `dataframe.describe()` method to view the percentiles.

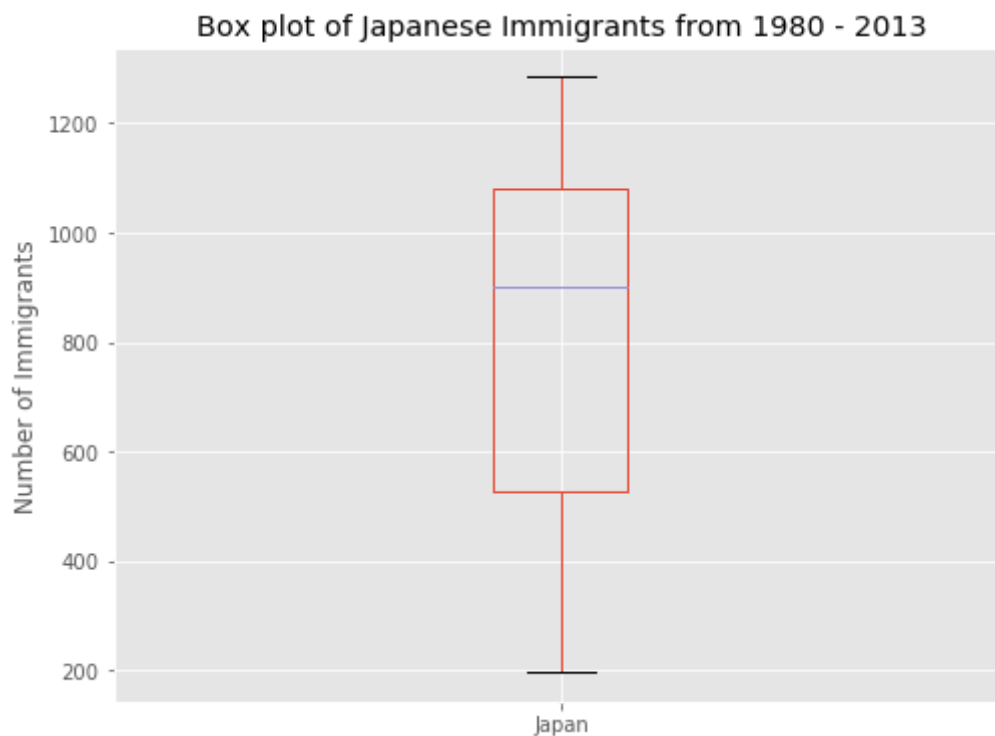
In [16]:

Out[16]:

Country	Japan
1980	701
1981	756
1982	598
1983	309
1984	246

Step 2: Plot by passing in `kind='box'` .

In [17]:



We can immediately make a few key observations from the plot above:

1. The minimum number of immigrants is around 200 (min), maximum number is around 1300 (max), and median number of immigrants is around 900 (median).
2. 25% of the years for period 1980 - 2013 had an annual immigrant count of ~500 or fewer (First quartile).
3. 75% of the years for period 1980 - 2013 had an annual immigrant count of ~1100 or fewer (Third quartile).

We can view the actual numbers by calling the `describe()` method on the dataframe.

In [18]:

Out[18]:

Country	Japan
count	34.000000
mean	814.911765
std	337.219771
min	198.000000
25%	529.000000
50%	902.000000
75%	1079.000000
max	1284.000000

One of the key benefits of box plots is comparing the distribution of multiple datasets. In one of the previous labs, we observed that China and India had very similar immigration trends. Let's analyze these two countries further using box plots.

Question: Compare the distribution of the number of new immigrants from India and China for the period 1980 - 2013.

Step 1: Get the dataset for China and India and call the dataframe **df_CI**.

In [19]:

Out[19]:

Country	China	India
1980	5123	8880
1981	6682	8670
1982	3308	8147
1983	1863	7338
1984	1527	5704

Double-click **here** for the solution.

Let's view the percentages associated with both countries using the `describe()` method.

In [20]:

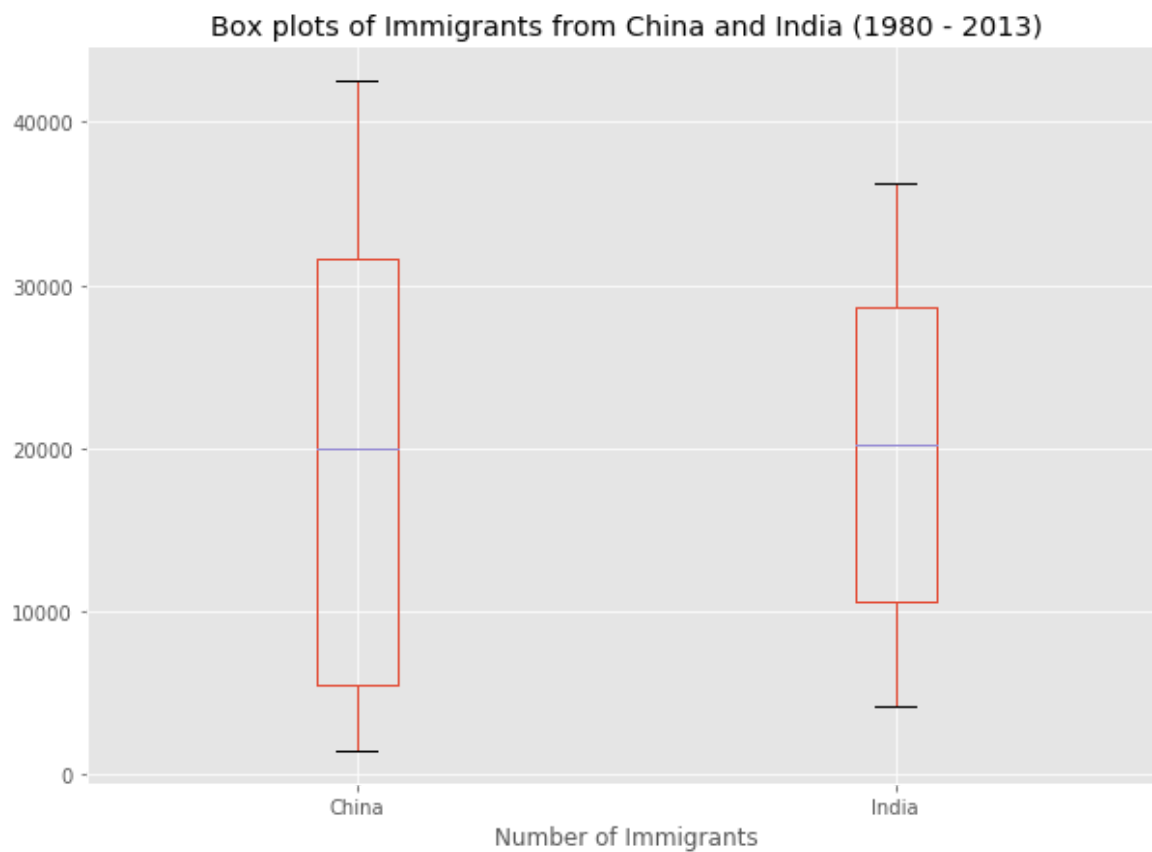
Out[20]:

Country	China	India
count	34.000000	34.000000
mean	19410.647059	20350.117647
std	13568.230790	10007.342579
min	1527.000000	4211.000000
25%	5512.750000	10637.750000
50%	19945.000000	20235.000000
75%	31568.500000	28699.500000
max	42584.000000	36210.000000

Double-click **here** for the solution.

Step 2: Plot data.

In [21]:

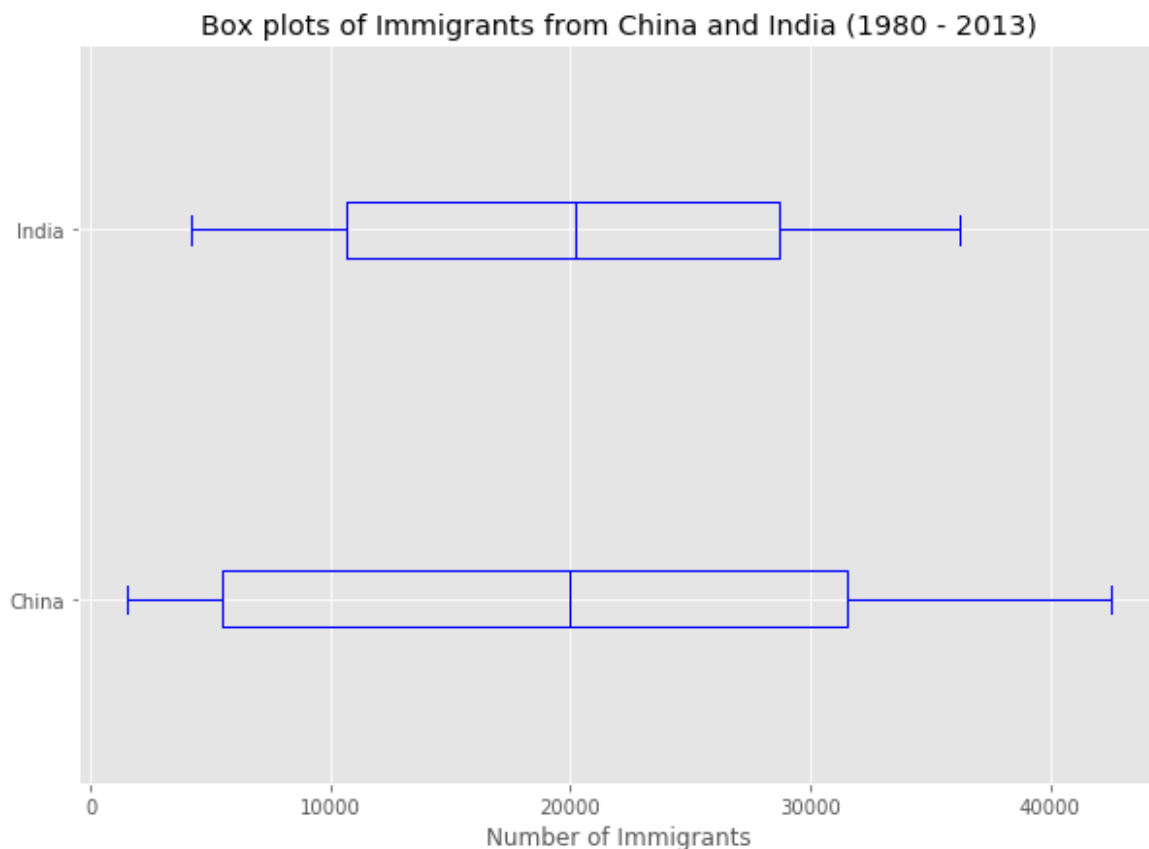


Double-click **here** for the solution.

We can observe that, while both countries have around the same median immigrant population (~20,000), China's immigrant population range is more spread out than India's. The maximum population from India for any year (36,210) is around 15% lower than the maximum population from China (42,584).

If you prefer to create horizontal box plots, you can pass the `vert` parameter in the **plot** function and assign it to *False*. You can also specify a different color in case you are not a big fan of the default red color.

In [22]:



Subplots

Often times we might want to plot multiple plots within the same figure. For example, we might want to perform a side by side comparison of the box plot with the line plot of China and India's immigration.

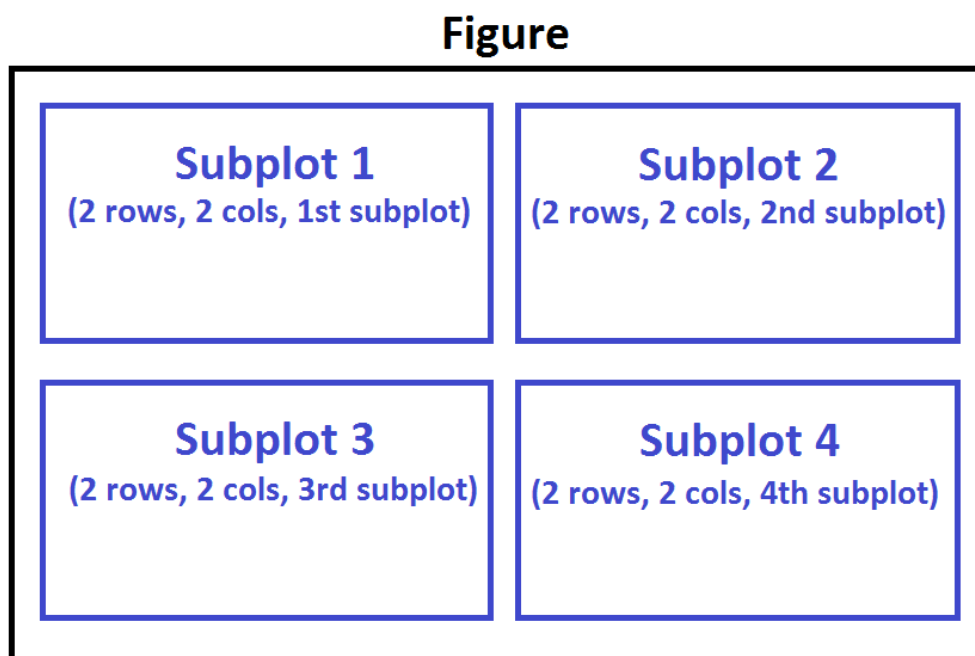
To visualize multiple plots together, we can create a **figure** (overall canvas) and divide it into **subplots**, each containing a plot. With **subplots**, we usually work with the **artist layer** instead of the **scripting layer**.

Typical syntax is :

```
fig = plt.figure() # create figure
ax = fig.add_subplot(nrows, ncols, plot_number) # create subplots
```

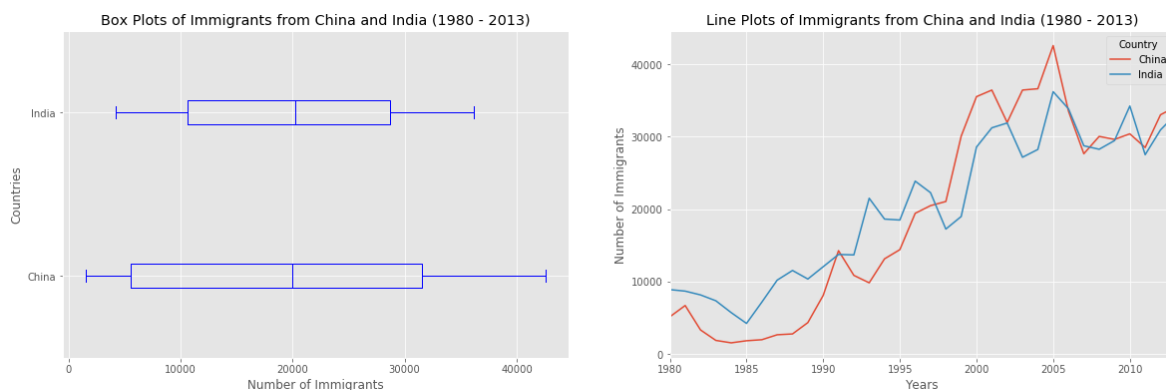
Where

- `nrows` and `ncols` are used to notionally split the figure into (`nrows * ncols`) sub-axes,
- `plot_number` is used to identify the particular subplot that this function is to create within the notional grid. `plot_number` starts at 1, increments across rows first and has a maximum of `nrows * ncols` as shown below.



We can then specify which subplot to place each plot by passing in the `ax` parameter in `plot()` method as follows:

In [23]:



* Tip regarding subplot convention

In the case when `nrows`, `ncols`, and `plot_number` are all less than 10, a convenience exists such that the a 3 digit number can be given instead, where the hundreds represent `nrows`, the tens represent `ncols` and the units represent `plot_number`. For instance,

```
subplot(211) == subplot(2, 1, 1)
```

produces a subaxes in a figure which represents the top plot (i.e. the first) in a 2 rows by 1 column notional grid (no grid actually exists, but conceptually this is how the returned subplot has been positioned).

Let's try something a little more advanced.

Previously we identified the top 15 countries based on total immigration from 1980 - 2013.

Question: Create a box plot to visualize the distribution of the top 15 countries (based on total immigration) grouped by the *decades* 1980s, 1990s, and 2000s.

Step 1: Get the dataset. Get the top 15 countries based on Total immigrant population. Name the dataframe **df_top15**.

In [24]:

Out[24]:

	Continent	Region	DevName	1980	1981	1982	1983	1984	1985	1986
Country										
India	Asia	Southern Asia	Developing regions	8880	8670	8147	7338	5704	4211	7150
China	Asia	Eastern Asia	Developing regions	5123	6682	3308	1863	1527	1816	1960
United Kingdom of Great Britain and Northern Ireland	Europe	Northern Europe	Developed regions	22045	24796	20620	10015	10170	9564	9470
Philippines	Asia	South-Eastern Asia	Developing regions	6051	5921	5249	4562	3801	3150	4166
Pakistan	Asia	Southern Asia	Developing regions	978	972	1201	900	668	514	691
United States of America	Northern America	Northern America	Developed regions	9378	10030	9074	7100	6661	6543	7074
Iran (Islamic Republic of)	Asia	Southern Asia	Developing regions	1172	1429	1822	1592	1977	1648	1794
Sri Lanka	Asia	Southern Asia	Developing regions	185	371	290	197	1086	845	1838
Republic of Korea	Asia	Eastern Asia	Developing regions	1011	1456	1572	1081	847	962	1208
Poland	Europe	Eastern Europe	Developed regions	863	2930	5881	4546	3588	2819	4808
Lebanon	Asia	Western Asia	Developing regions	1409	1119	1159	789	1253	1683	2576
France	Europe	Western Europe	Developed regions	1729	2027	2219	1490	1169	1177	1298
Jamaica	Latin America and the Caribbean	Caribbean	Developing regions	3198	2634	2661	2455	2508	2938	4649
Viet Nam	Asia	South-Eastern Asia	Developing regions	1191	1829	2162	3404	7583	5907	2741
Romania	Europe	Eastern Europe	Developed regions	375	438	583	543	524	604	656

15 rows × 38 columns

Double-click **here** for the solution.

Step 2: Create a new dataframe which contains the aggregate for each decade. One way to do that:

1. Create a list of all years in decades 80's, 90's, and 00's.
2. Slice the original dataframe `df_can` to create a series for each decade and sum across all years for each country.
3. Merge the three series into a new data frame. Call your dataframe **new_df**.

In [25]:

Out[25]:

	1980s	1990s	2000s
Country			
India	82154	180395	303591
China	32003	161528	340385
United Kingdom of Great Britain and Northern Ireland	179171	261966	83413
Philippines	60764	138482	172904
Pakistan	10591	65302	127598

Double-click **here** for the solution.

Let's learn more about the statistics associated with the dataframe using the `describe()` method.

In [26]:

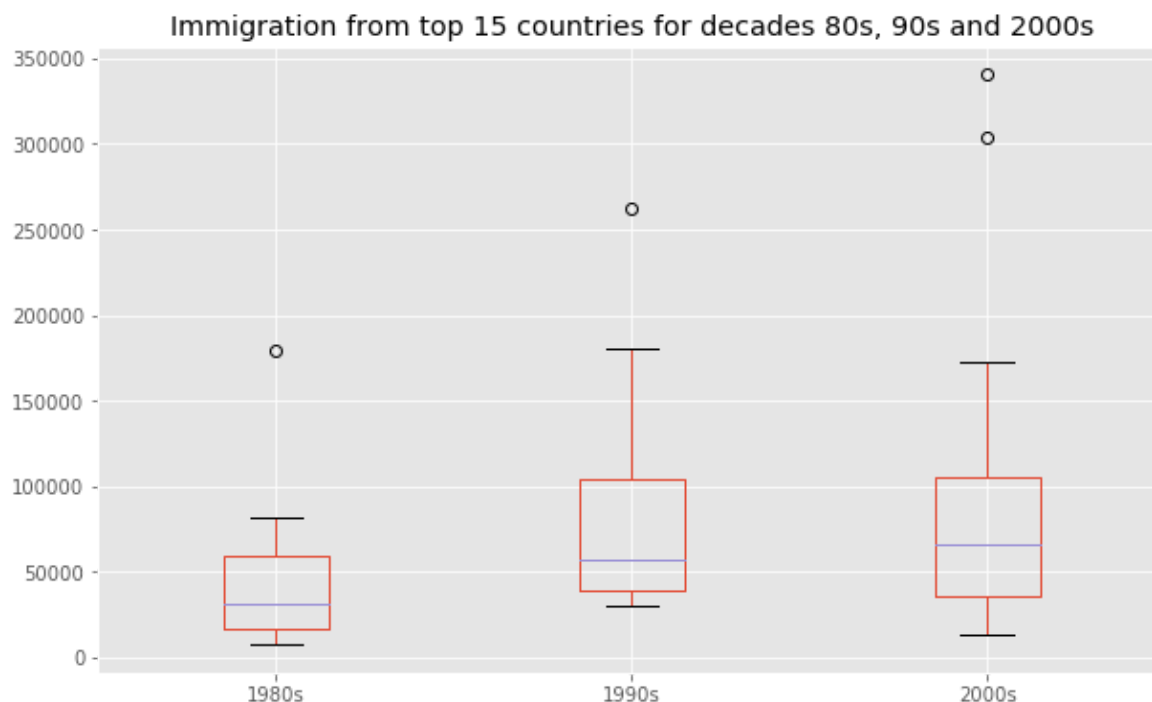
Out[26]:

	1980s	1990s	2000s
count	15.000000	15.000000	15.000000
mean	44418.333333	85594.666667	97471.533333
std	44190.676455	68237.560246	100583.204205
min	7613.000000	30028.000000	13629.000000
25%	16698.000000	39259.000000	36101.500000
50%	30638.000000	56915.000000	65794.000000
75%	59183.000000	104451.500000	105505.500000
max	179171.000000	261966.000000	340385.000000

Double-click **here** for the solution.

Step 3: Plot the box plots.

In [27]:



Double-click **here** for the solution.

Note how the box plot differs from the summary table created. The box plot scans the data and identifies the outliers. In order to be an outlier, the data value must be:

- larger than Q3 by at least 1.5 times the interquartile range (IQR), or,
- smaller than Q1 by at least 1.5 times the IQR.

Let's look at decade 2000s as an example:

- Q1 (25%) = 36,101.5
- Q3 (75%) = 105,505.5
- IQR = Q3 - Q1 = 69,404

Using the definition of outlier, any value that is greater than Q3 by 1.5 times IQR will be flagged as outlier.

Outlier > $105,505.5 + (1.5 * 69,404)$

Outlier > 209,611.5

In [28]:

Out[28]:

	1980s	1990s	2000s
Country			
India	82154	180395	303591
China	32003	161528	340385

China and India are both considered as outliers since their population for the decade exceeds 209,611.5.

The box plot is an advanced visualization tool, and there are many options and customizations that exceed the scope of this lab. Please refer to [Matplotlib documentation](http://matplotlib.org/api/pyplot_api.html#matplotlib.pyplot.boxplot) (http://matplotlib.org/api/pyplot_api.html#matplotlib.pyplot.boxplot) on box plots for more information.

Scatter Plots

A `scatter plot` (2D) is a useful method of comparing variables against each other. `Scatter plots` look similar to `line plots` in that they both map independent and dependent variables on a 2D graph. While the datapoints are connected together by a line in a line plot, they are not connected in a scatter plot. The data in a scatter plot is considered to express a trend. With further analysis using tools like regression, we can mathematically calculate this relationship and use it to predict trends outside the dataset.

Let's start by exploring the following:

Using a `scatter plot`, let's visualize the trend of total immigration to Canada (all countries combined) for the years 1980 - 2013.

Step 1: Get the dataset. Since we are expecting to use the relationship between `years` and `total population`, we will convert `years` to `int` type.

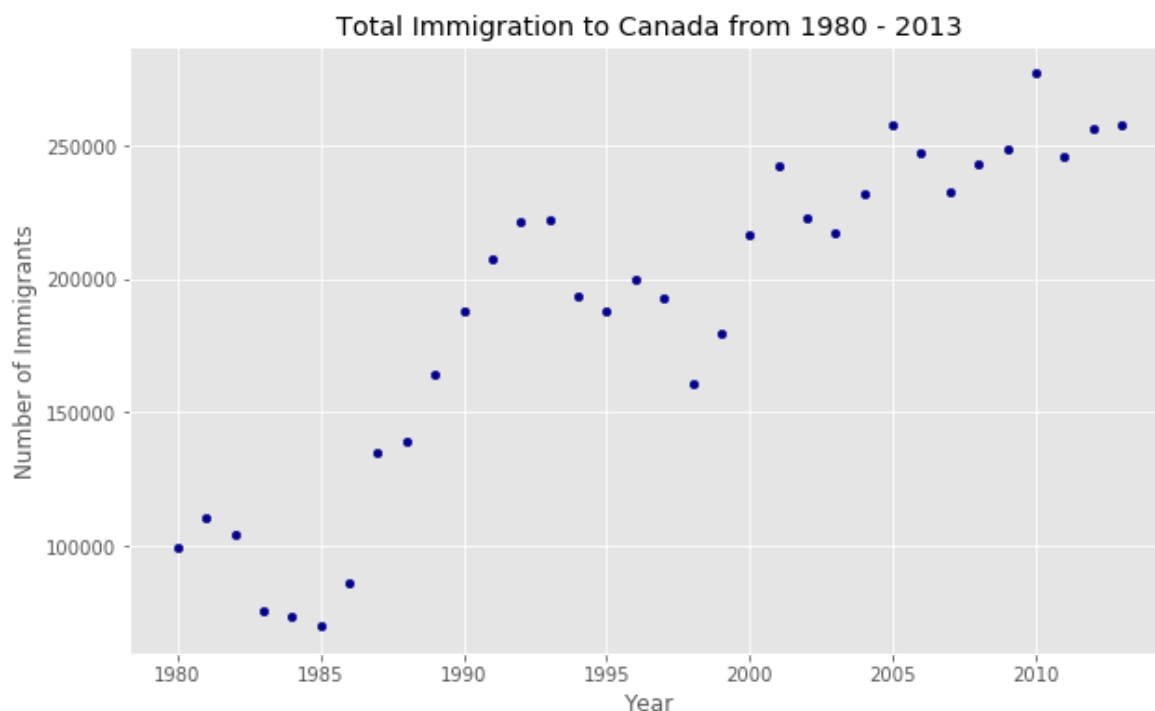
In [29]:

Out[29]:

	year	total
0	1980	99137
1	1981	110563
2	1982	104271
3	1983	75550
4	1984	73417

Step 2: Plot the data. In `Matplotlib`, we can create a `scatter` plot set by passing in `kind='scatter'` as plot argument. We will also need to pass in `x` and `y` keywords to specify the columns that go on the x- and the y-axis.

In [30]:



Notice how the scatter plot does not connect the datapoints together. We can clearly observe an upward trend in the data: as the years go by, the total number of immigrants increases. We can mathematically analyze this upward trend using a regression line (line of best fit).

So let's try to plot a linear line of best fit, and use it to predict the number of immigrants in 2015.

Step 1: Get the equation of line of best fit. We will use **Numpy's** `polyfit()` method by passing in the following:

- `x` : x-coordinates of the data.
- `y` : y-coordinates of the data.
- `deg` : Degree of fitting polynomial. 1 = linear, 2 = quadratic, and so on.

In [31]:

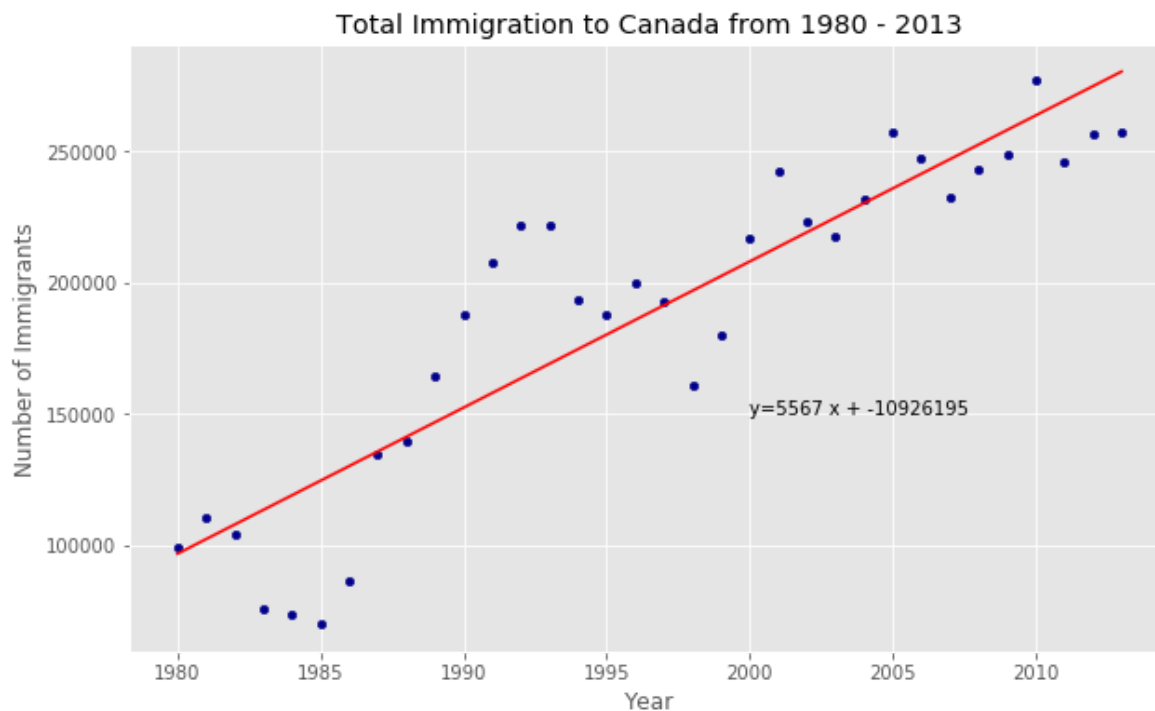
Out[31]:

```
array([ 5.56709228e+03, -1.09261952e+07])
```

The output is an array with the polynomial coefficients, highest powers first. Since we are plotting a linear regression $y = a \cdot x + b$, our output has 2 elements $[5.56709228e+03, -1.09261952e+07]$ with the slope in position 0 and intercept in position 1.

Step 2: Plot the regression line on the scatter plot .

In [32]:



Out[32]:

```
'No. Immigrants = 5567 * Year + -10926195'
```

Using the equation of line of best fit, we can estimate the number of immigrants in 2015:

```
No. Immigrants = 5567 * Year - 10926195
No. Immigrants = 5567 * 2015 - 10926195
No. Immigrants = 291,310
```

When compared to the actuals from Citizenship and Immigration Canada's (CIC) [2016 Annual Report](http://www.cic.gc.ca/english/resources/publications/annual-report-2016/index.asp) (<http://www.cic.gc.ca/english/resources/publications/annual-report-2016/index.asp>), we see that Canada accepted 271,845 immigrants in 2015. Our estimated value of 291,310 is within 7% of the actual number, which is pretty good considering our original data came from United Nations (and might differ slightly from CIC data).

As a side note, we can observe that immigration took a dip around 1993 - 1997. Further analysis into the topic revealed that in 1993 Canada introduced Bill C-86 which introduced revisions to the refugee determination system, mostly restrictive. Further amendments to the Immigration Regulations cancelled the sponsorship required for "assisted relatives" and reduced the points awarded to them, making it more difficult for family members (other than nuclear family) to immigrate to Canada. These restrictive measures had a direct impact on the immigration numbers for the next several years.

Question: Create a scatter plot of the total immigration from Denmark, Norway, and Sweden to Canada from 1980 to 2013?

Step 1: Get the data:

1. Create a dataframe the consists of the numbers associated with Denmark, Norway, and Sweden only. Name it **df_countries**.
2. Sum the immigration numbers across all three countries for each year and turn the result into a dataframe. Name this new dataframe **df_total**.
3. Reset the index in place.
4. Rename the columns to **year** and **total**.
5. Display the resulting dataframe.

In [33]:

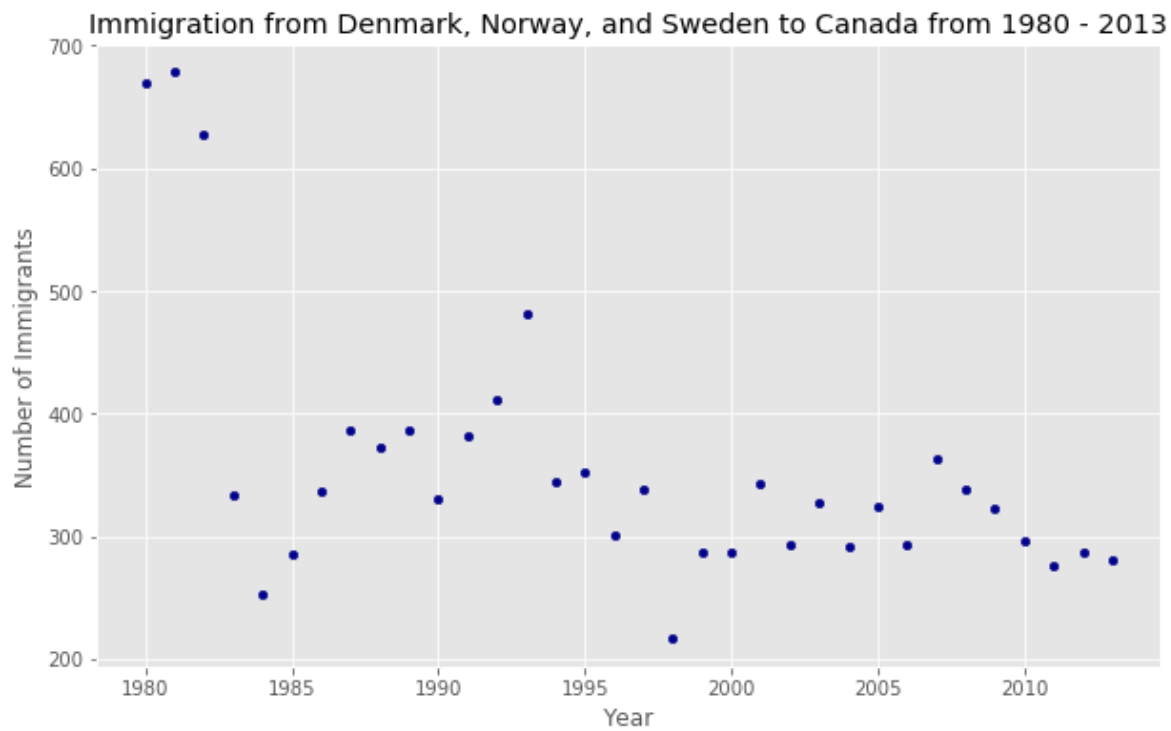
Out[33]:

	year	total
0	1980	669
1	1981	678
2	1982	627
3	1983	333
4	1984	252

Double-click **here** for the solution.

Step 2: Generate the scatter plot by plotting the total versus year in **df_total**.

In [34]:



Double-click **here** for the solution.

Bubble Plots

A `bubble plot` is a variation of the `scatter plot` that displays three dimensions of data (x, y, z). The datapoints are replaced with bubbles, and the size of the bubble is determined by the third variable 'z', also known as the weight. In `matplotlib`, we can pass in an array or scalar to the keyword `s` to `plot()`, that contains the weight of each point.

Let's start by analyzing the effect of Argentina's great depression.

Argentina suffered a great depression from 1998 - 2002, which caused widespread unemployment, riots, the fall of the government, and a default on the country's foreign debt. In terms of income, over 50% of Argentines were poor, and seven out of ten Argentine children were poor at the depth of the crisis in 2002.

Let's analyze the effect of this crisis, and compare Argentina's immigration to that of it's neighbour Brazil. Let's do that using a `bubble plot` of immigration from Brazil and Argentina for the years 1980 - 2013. We will set the weights for the bubble as the *normalized* value of the population for each year.

Step 1: Get the data for Brazil and Argentina. Like in the previous example, we will convert the `years` to type `int` and bring it in the dataframe.

In [35]:

Out[35]:

Country	Year	Afghanistan	Albania	Algeria	American Samoa	Andorra	Angola	Antigua and Barbuda	Argentina
0	1980	16	1	80	0	0	1	0	368
1	1981	39	0	67	1	0	3	0	426
2	1982	39	0	71	0	0	6	0	626
3	1983	47	0	69	0	0	6	0	241
4	1984	71	0	63	0	0	4	42	237

5 rows × 10 columns

Step 2: Create the normalized weights.

There are several methods of normalizations in statistics, each with its own use. In this case, we will use [feature scaling](https://en.wikipedia.org/wiki/Feature_scaling) (https://en.wikipedia.org/wiki/Feature_scaling) to bring all values into the range [0,1]. The general formula is:

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

where x is an original value, x' is the normalized value. The formula sets the max value in the dataset to 1, and sets the min value to 0. The rest of the datapoints are scaled to a value between 0-1 accordingly.

In [36]:

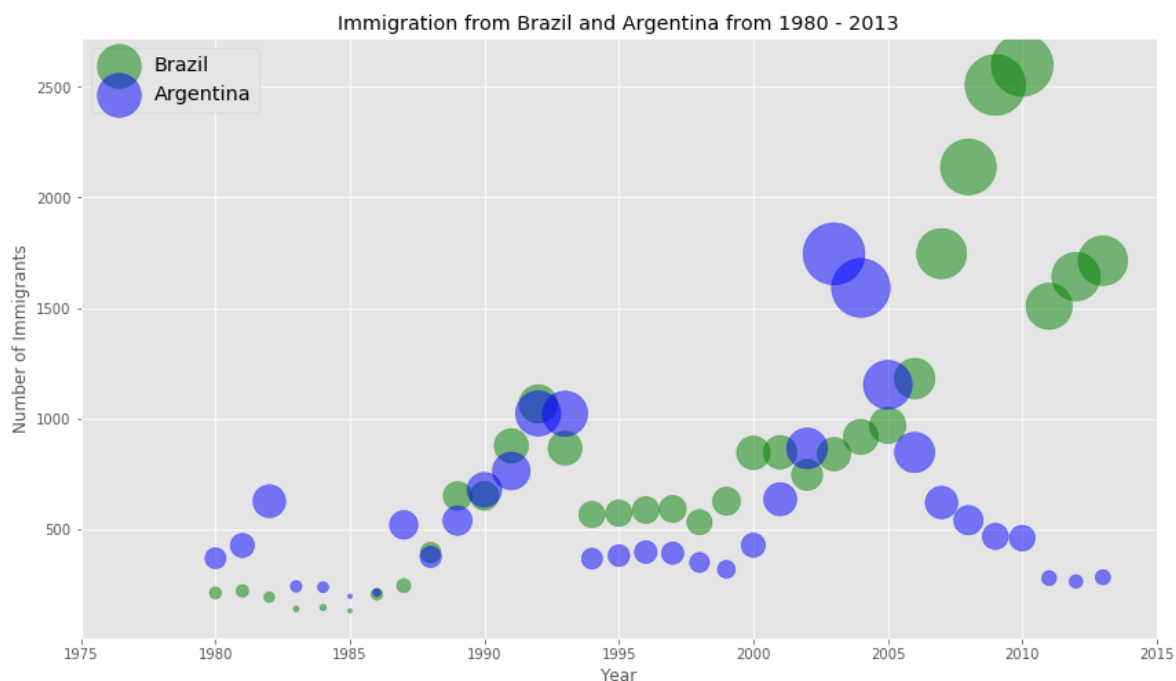
Step 3: Plot the data.

- To plot two different scatter plots in one plot, we can include the axes one plot into the other by passing it via the `ax` parameter.
- We will also pass in the weights using the `s` parameter. Given that the normalized weights are between 0-1, they won't be visible on the plot. Therefore we will:
 - multiply weights by 2000 to scale it up on the graph, and,
 - add 10 to compensate for the min value (which has a 0 weight and therefore scale with x2000).

In [37]:

Out[37]:

<matplotlib.legend.Legend at 0x7fa4485fb208>



The size of the bubble corresponds to the magnitude of immigrating population for that year, compared to the 1980 - 2013 data. The larger the bubble, the more immigrants in that year.

From the plot above, we can see a corresponding increase in immigration from Argentina during the 1998 - 2002 great depression. We can also observe a similar spike around 1985 to 1993. In fact, Argentina had suffered a great depression from 1974 - 1990, just before the onset of 1998 - 2002 great depression.

On a similar note, Brazil suffered the *Samba Effect* where the Brazilian real (currency) dropped nearly 35% in 1999. There was a fear of a South American financial crisis as many South American countries were heavily dependent on industrial exports from Brazil. The Brazilian government subsequently adopted an austerity program, and the economy slowly recovered over the years, culminating in a surge in 2010. The immigration data reflect these events.

Question: Previously in this lab, we created box plots to compare immigration from China and India to Canada. Create bubble plots of immigration from China and India to visualize any differences with time from 1980 to 2013. You can use `df_can_t` that we defined and used in the previous example.

Step 1: Normalize the data pertaining to China and India.

In [38]:

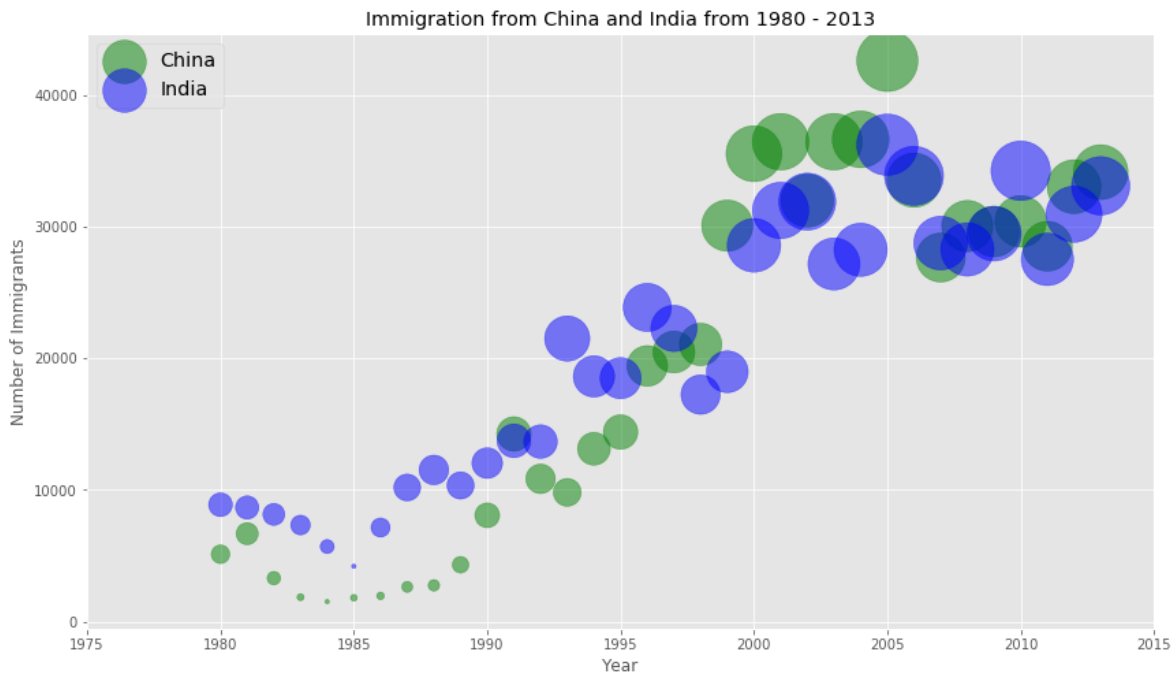
Double-click **here** for the solution.

Step 2: Generate the bubble plots.

In [39]:

Out[39]:

<matplotlib.legend.Legend at 0x7fa44895df28>



Double-click **here** for the solution.

Thank you for completing this lab!

This notebook was created by [Jay Rajasekharan](https://www.linkedin.com/in/jayrajasekharan) (<https://www.linkedin.com/in/jayrajasekharan>) with contributions from [Ehsan M. Kermani](https://www.linkedin.com/in/ehsanmkermani) (<https://www.linkedin.com/in/ehsanmkermani>), and [Slobodan Markovic](https://www.linkedin.com/in/slobodan-markovic) (<https://www.linkedin.com/in/slobodan-markovic>).

This notebook was recently revamped by [Alex Aklson](https://www.linkedin.com/in/aklson/) (<https://www.linkedin.com/in/aklson/>). I hope you found this lab session interesting. Feel free to contact me if you have any questions!

This notebook is part of a course on **Coursera** called *Data Visualization with Python*. If you accessed this notebook outside the course, you can take this course online by clicking [here](http://cocl.us/DV0101EN_Coursera_Week2_LAB2) (http://cocl.us/DV0101EN_Coursera_Week2_LAB2).

Copyright © 2019 [Cognitive Class](https://cognitiveclass.ai/?utm_source=bducopyrightlink&utm_medium=dswb&utm_campaign=bdu) (https://cognitiveclass.ai/?utm_source=bducopyrightlink&utm_medium=dswb&utm_campaign=bdu). This notebook and its source code are released under the terms of the [MIT License](https://bigdatauniversity.com/mit-license/) (<https://bigdatauniversity.com/mit-license/>).