

Data visualisation in Python

Prof. Dorien Herremans

For a full overview of types of plots using matplotlib, see the gallery at <https://matplotlib.org/2.0.2/gallery.html>

Scatterplots

We will be using scottish_hills.csv from <https://github.com/ourcodingclub/CC-python-pandas-matplotlib>. The file contains all the mountains above 3000 feet (about 914 metres) in Scotland.

We can read this into a variable and see the first 10 lines:

```
In [1]: 1 import pandas as pd
2         dataframe = pd.read_csv("scottish_hills.csv")
3         print(dataframe.head(10))

0      Hill Name  Height  Latitude  Longitude  Osgrid
1      A' Bhuidheanach Beag  936.0  56.870362  -4.199001  NM660775
2      A' Chailleach  997.0  57.693800  -5.128715  NM136714
3      A' Chailleach  929.2  57.109564  -4.179285  NM681041
4      A' Chailleach  1120.0  57.184186  -5.154837  NM094167
5      A' Ghlas-bheinn  918.0  57.255090  -5.303687  NM080231
6      A' Mhallaighean  967.0  57.718664  -5.146720  NM077449
7      A' Mharcaonach  973.2  56.857002  -4.290668  NM664762
8      Am Basteir  936.0  57.247931  -6.202962  NM665253
9      Am Bodach  1031.8  56.741727  -4.983393  NM176650
10     Am Faochagach  953.0  57.771801  -4.853899  NM303793
```

As explored last week, pandas dataframes can be used for some preliminary data exploration. For instance, let's sort the hills by height:

```
In [2]: 1 sorted_hills = dataframe.sort_values(by=['Height'], ascending=False)
2         # Let's have a look at the top 5 to check
3         print(sorted_hills.head(5))

      Hill Name  Height  Latitude  Longitude  Osgrid
92  Ben Nevis  1344.5  56.796891  -5.003675  NM166712
88  Ben Macdui (Beninn Macduibh)  1309.0  57.070368  -3.669039  NM988989
104 Brazerach  1296.0  57.078298  -3.728389  NM953999
212 Cairn Doui  1291.0  57.054397  -3.710773  NM963972
212 Sgor an Lochain Uaine  1256.0  57.059369  -3.725797  NM954976
```

Now let's load matplotlib. Note: if you are using a jupyter notebook you need the inline statement on line 1 below:

```
In [3]: 1 %matplotlib inline
2         import matplotlib.pyplot as plt

In [4]: 1 x = dataframe.Height
2         y = dataframe.Latitude
3         z = dataframe.Longitude
```

Now we are ready to start visualising them. Let's create (and save) a scatterplot:

```
In [5]: 1 plt.scatter(x, y)
2         plt.savefig("scottish_scatter_plot.png")

58.5
58.0
57.5
57.0
56.5

900 1000 1100 1200 1300
```

If you are not using Python, you can use `plt.show()` to display the plot.

Now let's build upon this graph by adding a linear regression line to it.

```
In [6]: 1 from scipy.stats import linregress
2         stats = linregress(x, y)
3         # Add x and y labels, and set their font size
4         m = stats.slope
5         b = stats.intercept

Now we can add the plot of our linear regression by using the equation of a straight line:
```

```
In [7]: 1 plt.scatter(x, y)
2         plt.plot(x, m * x + b, color='red') # The equation of the straight line.

Out[7]: <matplotlib.lines.Line2D at 0x1102d1f98>
```

Note, whether this line is statistically significant can be determined using the extra information in the stats object - `stats.pvalue` and `stats.rvalue`.

Now you can make your plot look nicer using arguments such as `fontsize`, `linewidth`, `color`...

```
In [8]: 1 # Change the default figure size
2         plt.figure(figsize=(10,10))
3         # Change the default marker for the scatter from circles to 'x's
4         plt.scatter(x, y, marker='x')
5         # Set the linewidth on the regression line to 3px
6         plt.plot(x, m * x + b, color='red', linewidth=3)
7         # Add x and y labels, and set their font size
8         plt.xlabel("Height (m)", fontsize=20)
9         plt.ylabel("Latitude", fontsize=20)
10        # Create the new height variable
11        # Creating the new height variable
12        # Setting the font size of the number labels on the axes
13        # Set the font size of the number labels on the axes
14        # Set the font size of the number labels on the axes
15        plt.xticks(fontsize=18)
16        plt.yticks(fontsize=18)

Out[8]: (array([ 56. ,  56.5,  57. ,  57.5,  58. ,  58.5,  59. ]),
<list of 7 TextYtickLabel objects>)
```

Let's have a look at how the hills are spread out geographically using latitude (y) and longitude (z). Now we can use `s = x` to say that the size needs to be equal to the height (x), (8 added -900 to make the difference between big and small mountains larger)

```
In [9]: 1 import numpy as np
2         colors = np.random.rand(len(y)) # generates a different color for each different mountain
3         plt.scatter(y, z, s = (x-900), c=colors, alpha=0.5)

Out[9]: <matplotlib.collections.PathCollection at 0x110581a20>
```

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

```
In [10]: 1 plt.hist(x, bins=25, normed=True) # bins separates the latitude in 25 discrete categories. Normed will normalize the
2         plt.savefig("histogram.png", dpi=25) # results in 160x120 px image

0.007
0.006
0.005
0.004
0.003
0.002
0.001
0.000

900 1000 1100 1200 1300
```

Quickly style your plot with stylesheets, full overview at https://matplotlib.org/gallery/style_sheets/style_sheets_reference.html.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

```
In [11]: 1 import numpy as np
2         # using a stylesheet
3         plt.style.use('seaborn-pastel')
4         # Creating the new height variable
5         # Setting the font size of the number labels on the axes
6         # Setting the font size of the number labels on the axes
7         # Setting the font size of the number labels on the axes
8         fig, ax = plt.subplots()
9         ax.hist(x, bins=25, normed=True, histtype='stepfilled', alpha=0.8, label='Height')
10        ax.hist(shifted_x, bins=25, normed=True, histtype='stepfilled', alpha=0.8, label='Height - 100')
11        ax.legend(prop={'size': 10})
12        # Setting the font size of the number labels on the axes
13        # Setting the font size of the number labels on the axes
14        ax.set_xlabel('Normalized distribution')
15        ax.set_ylabel('Height')

Out[11]: <matplotlib.text.Text at 0x11044f8d0>
```

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.

Let's also create a new variable that contains the height of the hills - 100m. This to illustrate how to add a second distribution to your graph. In this case, we will make them slightly transparent.