

数据可视化

作业二

姓名: 冉诗菡

学号: 15307130424

计算机科学与技术 (数据科学方向)

复旦大学

大数据学院

2018 年 4 月 17 日

题目一：熟悉Matplotlib作图

题目描述

Find/design 5 sets of different data, and use 5 different types of plots to visualize the data using Python and matplotlib; please take a few sentences to describe the data information, background, and visualization effects for analysis. Submit your 5 data sets and code.

解答

首先我们来看一下Matplotlib作图的基本流程：

```
1  from matplotlib import pyplot as plt
2
3  #创建坐标轴
4  fig, ax = plt.subplots()
5
6  #1.显示整个窗口的标题
7  fig.canvas.set_window_title("window_title")
8
9  #2.显示当前图表的x, y轴标注
10 ax.set_xlabel("xlabel")
11 ax.set_ylabel("ylabel")
12
13 #3.显示坐标轴的网格
14 ax.grid()
15
16 #set title of the figure
17 ax.set_title("figure_title")
18
19 #从数据源中的那个国家的开始点和结束点开始统计
20 start, end= 10, 20
21
22 #4.标注所有线的标记类型，方便不便区分颜色的人阅读
23 markerSymbols = ['.', ',', 'o', 'v', '^', '<', '>', 's', 'p', 'P']
24
25 #5.设置x轴的大小范围
26 ax.set_xlim(min(x_pos), max(x_pos))
27
28 #6.目前颜色自动分配
29
30 #开始画图
```

```

31 for i,gdps in enumerate(pdata.iloc[startCountry:endCountry].values):
32     line, = ax.plot([], [], lw=1) # 感觉最常用 plt.plot
33     #给每条线标注线的标记类型
34     line.set_marker(markerSymbols[i])
35     #给每条线填充data
36     line.set_data(x_pos[0:], gdps[from_:])
37
38 #7.显示每种线的缩略图标记
39 #linesets是各个数据集画出的不同线，countries是每条线对应的国家名字
40 plt.legend(linesets, countries, loc='upper left')
41
42 #8.更改数据的最小最大值范围，使其呈居中水平
43 delta = maxY - minY
44 minY = minY - delta / 2
45 maxY = maxY + delta
46
47 # set_ylim --- 没有这个会很有问题
48 ax.set_ylim(minY, maxY)
49
50 #9.展示图片
51 plt.show()

```

作图一：散点图

生成1024个呈标准正态分布的二维数据组(平均数是0，方差为1)作为一个数据集。

完善：C = np.arctan2(Y,X)保证了散点图可以根据散点所在的位置选择点的颜色。

```

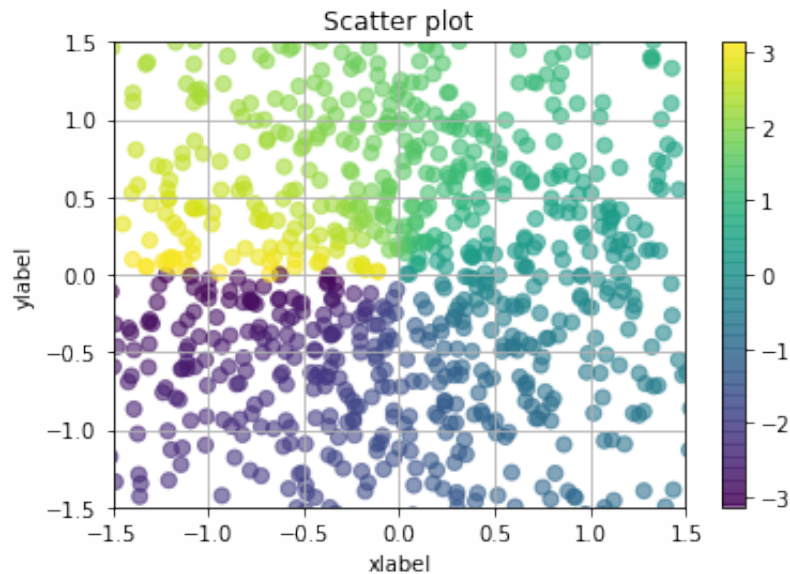
1  import matplotlib.pyplot as plt
2  import numpy as np
3
4  ##### 数据生成 #####
5  n = 1024 # data size
6  X = np.random.normal(0, 1, n) # 每一个点的X值
7  Y = np.random.normal(0, 1, n) # 每一个点的Y值
8  C = np.arctan2(Y, X) # for color value
9
10 ##### 作图过程 #####
11 plt.scatter(X, Y, s=50, c=C, alpha=.6)
12
13 plt.xlim(-1.5, 1.5)
14 plt.ylim(-1.5, 1.5)
15
16 plt.xlabel("xlabel")
17 plt.ylabel("ylabel")
18
19 plt.grid()
20 plt.colorbar()

```

```

21 plt.title("Scatter plot")
22
23 plt.show()

```



作图二：柱状图

模拟月度报表，向上向下分别生成12个数据，X为0到11的整数，Y是相应的均匀分布的随机数据。

完善：给柱状图更换颜色、网格、图例，并加上数字标注。

```

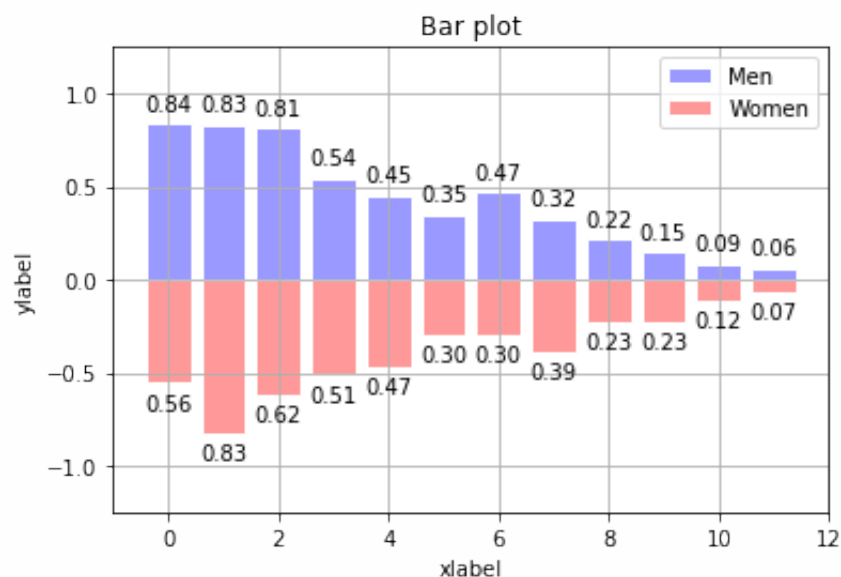
1  import matplotlib.pyplot as plt
2  import numpy as np
3
4  ##### 数据生成 #####
5  n = 12
6  X = np.arange(n)
7  Y1 = (1 - X / float(n)) * np.random.uniform(0.5, 1.0, n)
8  Y2 = (1 - X / float(n)) * np.random.uniform(0.5, 1.0, n)
9
10 ##### 作图过程 #####
11 # 用facecolor设置主体颜色，edgecolor设置边框颜色为白色
12 p1 = plt.bar(X, +Y1, facecolor='#9999ff', edgecolor='white')
13 p2 = plt.bar(X, -Y2, facecolor='#ff9999', edgecolor='white')
14
15 # 用函数plt.text分别在柱体上方（下方）加上数值
16 # 用%.2f保留两位小数，横向居中对齐ha='center'
17 # 纵向底部（顶部）对齐va='bottom'
18 for x, y in zip(X, Y1):
19     # ha: horizontal alignment
20     # va: vertical alignment
21     plt.text(x, y + 0.05, '%.2f' % y, ha='center', va='bottom')

```

```

22
23 for x, y in zip(X, Y2):
24     # ha: horizontal alignment
25     # va: vertical alignment
26     plt.text(x, -y - 0.05, '%.2f' % y, ha='center', va='top')
27
28 plt.xlim(-1, n)
29 plt.ylim(-1.25, 1.25)
30
31 plt.xlabel("xlabel")
32 plt.ylabel("ylabel")
33
34 plt.grid()
35 plt.title("Bar plot")
36 plt.legend((p1, p2), ('Men', 'Women'), loc='upper right')
37
38 plt.show()

```



作图三：饼图

这里用了课上Lab的一个例子，饼图是直观显示分布比例的一种可视化方式。

完善：给饼图添加颜色、标注、阴影、比例。通过explode参数把Logs呈现爆炸式的做法，给人造成了Logs和Frogs的比例差不多的视觉错觉。

```

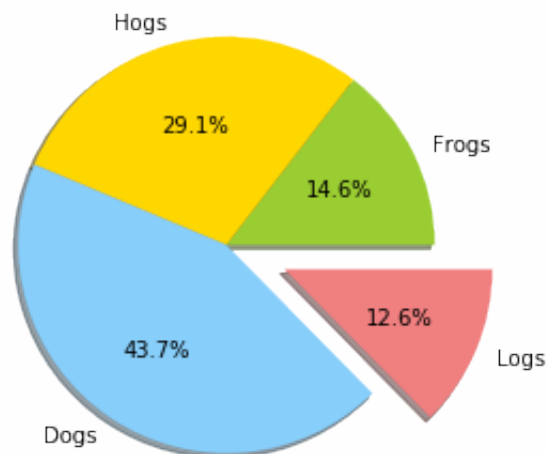
1 import matplotlib.pyplot as plt
2
3 ##### 数据生成 #####
4 labels = 'Frogs', 'Hogs', 'Dogs', 'Logs' # label的名称
5 sizes = [150, 300, 450, 130] # 比例数组

```

```

6 colors = ['yellowgreen', 'gold', 'lightskyblue', 'lightcoral'] # 颜色数组
7 explode = (0, 0, 0, 0.3) # 设定扇形图与中心点的距离
8
9 ##### 作图过程 #####
10 fig = plt.figure()
11
12 # 根据相关参数画出饼图
13 # matplotlib.pyplot.pie(x, explode=None, labels=None, colors=None, autopct=None, pctdistance=0.6,
14 # shadow=False, labeldistance=1.1, startangle=None, radius=None, counterclock=True, wedgeprops=None,
15 # textprops=None, center=(0, 0), frame=False, hold=None, data=None)
16
17 plt.pie(sizes, explode, labels, colors, shadow=True, autopct='%.1f%%')
18
19 # Set aspect ratio to be equal so that pie is drawn as a circle.
20 plt.axis('equal') # 显示方式, 如果为equql,则显示一个正圆,否则为椭圆
21
22 plt.show()

```



作图四：等高线图

数据集即三维点 (x,y) 和对应的高度值，共有256个点。高度值使用一个 height function $f(x,y)$ 生成。x, y 分别是在区间[-3,3]中均匀分布的256个值，并用meshgrid在二维平面中将每一个x和每一个y分别对应起来，编织成栅格。

完善：等高线进行颜色填充，增加线条和高度数值，省略X轴Y轴坐标。

```

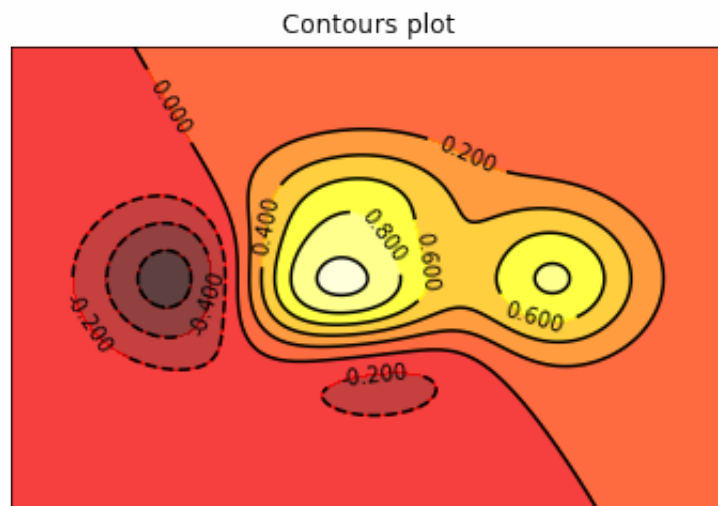
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4

```

```

5 ##### 数据生成 #####
6 def f(x,y):
7     # the height function
8     return (1 - x / 2 + x**5 + y**3) * np.exp(-x**2 -y**2)
9
10 n = 256
11 x = np.linspace(-3, 3, n)
12 y = np.linspace(-3, 3, n)
13 X,Y = np.meshgrid(x, y)
14
15 ##### 作图过程 #####
16 # 使用函数plt.contourf把颜色加进去, 位置参数分别为: X, Y, f(X,Y)
17 # 透明度0.75, 并将 f(X,Y) 的值对应到color map的暖色组中寻找对应颜色
18 plt.contourf(X, Y, f(X, Y), 8, alpha=.75, cmap=plt.cm.hot)
19
20 # 使用plt.contour函数进行等高线绘制
21 C = plt.contour(X, Y, f(X, Y), 8, colors='black')
22
23 plt.clabel(C, inline=True, fontsize=10)
24 plt.xticks(())
25 plt.yticks(())
26
27 plt.title("Contours plot")
28
29 plt.show()

```



作图五：3D数据

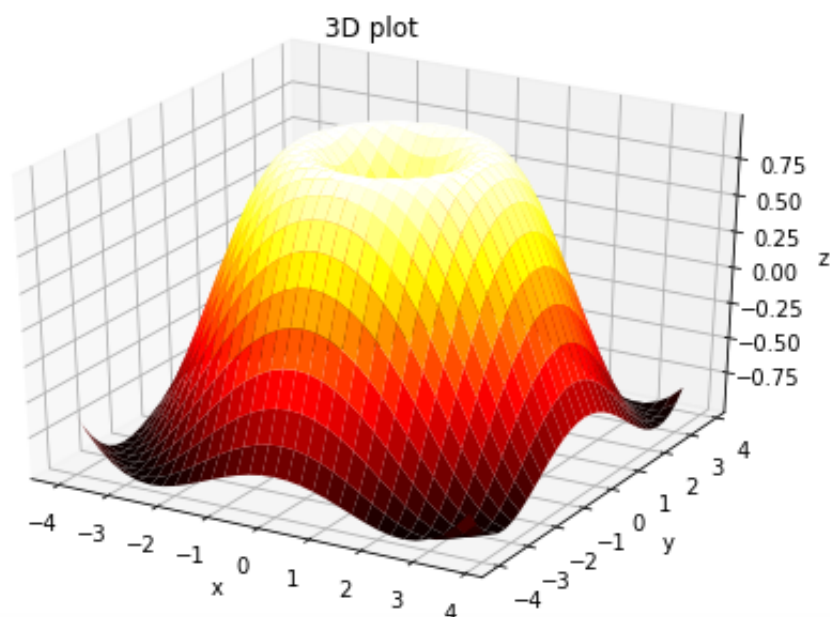
数据集即三维点 (x,y) 和对应的高度值，共有256个点。高度值使用一个 height function $f(x,y)$ 生成。 x, y 分别是在区间 $[-3,3]$ 中均匀分布的256个值，并用meshgrid在二维平面中将每一个x和每一个y分别对应起来，编织成栅格。

完善：等高线进行颜色填充，增加线条和高度数值，省略X轴Y轴坐标。

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from mpl_toolkits.mplot3d import Axes3D
4
5 ##### 数据生成 #####
6 # X, Y value
7 X = np.arange(-4, 4, 0.25)
8 Y = np.arange(-4, 4, 0.25)
9 # 将 X 和 Y 编织成栅格 即 x-y 平面的网格
10 X, Y = np.meshgrid(X, Y)
11 R = np.sqrt(X ** 2 + Y ** 2)
12 # 计算3D点的高度
13 Z = np.sin(R)
14
15 ##### 作图过程 #####
16 # 定义图像窗口, 在窗口上添加3D坐标轴
17 fig = plt.figure()
18 ax = Axes3D(fig)
19
20 # 做出一个三维曲面, 并用 colormap rainbow 填充颜色
21 # rstride 和 cstride 分别代表 row 和 column 的跨度
22 ax.plot_surface(X, Y, Z, rstride=1, cstride=1, cmap=plt.cm.hot)
23
24 # 设置图形展示效果
25 ax.set_xlabel('x')
26 ax.set_ylabel('y')
27 ax.set_zlabel('z')
28
29 ax.set_title('3D plot')
30
31 plt.show()

```



参考文献

- 0. [Matplotlib官方文档](#)
- 1. [莫烦Python - Matplotlib教程](#)

题目二：动画和隐喻的应用

题目描述

Find/design a dataset and visualize the data using either the techniques of animation or metaphor, or both of them. Please take a few sentences to describe the data information, background, and visualization effects for analysis. Submit your data and code. （做一个画图用到动画功能或（和）隐喻可视化功能）。

解答

代码如下，保存的动画在作业文件夹中。

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from matplotlib.animation import FuncAnimation
4
5 ##### 数据生成 #####
6 fig, ax = plt.subplots()
7 A = range(-5,5)
8 xdata, ydata = [], []
9 xdata2, ydata2 = [], []
10 xdata3, ydata3 = [], []
11
12 ln1, ln2, ln3, ln4 = ax.plot([], [], 'r-',
13                               [], [], 'b-',
14                               [], [], 'y-',
15                               [], [], 'c-',
16                               animated=False) #animated is associated with
17 blit
18
19 def init():
20     ax.set_xlim(0, 10)
21     ax.set_ylim(-6, 6)
22     return ln1, ln2, ln3, ln4
23
24 def update(i): #i is an int from 0 to frames-1, and keep looping
```

```

24     ax.set_xlim(i/250, 10+i/250)
25     global A
26     iter = int(i/50)
27
28     if i==0:
29         xdata.clear()
30         ydata.clear()
31         xdata2.clear()
32         ydata2.clear()
33         iter = 0
34
35     xdata.append(i/50)
36     ydata.append(A[iter])
37     xdata2.append(i/50)
38     ydata2.append(np.cos(i/100))
39     xdata3.append(i/50)
40     ydata3.append(np.cos(i/100-np.pi))
41     ln1.set_data(xdata, ydata)
42     ln2.set_data(xdata2, ydata2)
43     ln3.set_data(xdata3, ydata3)
44     x = np.linspace(0, 10, 1000)
45     y = np.sin(np.pi*(x + i/100))
46     ln4.set_data(x, y)
47     iter += 1
48     return ln1, ln2, ln3, ln4
49
50
51 ##### 作图过程 #####
52 def main():
53     ani = FuncAnimation(fig, update, frames = 500, interval = 20,
54                         init_func=init, blit=False)
55     ani.save('animation.mp4', writer='ffmpeg', fps=30)
56     plt.show()
57
58 if __name__ == '__main__':
59     main()

```

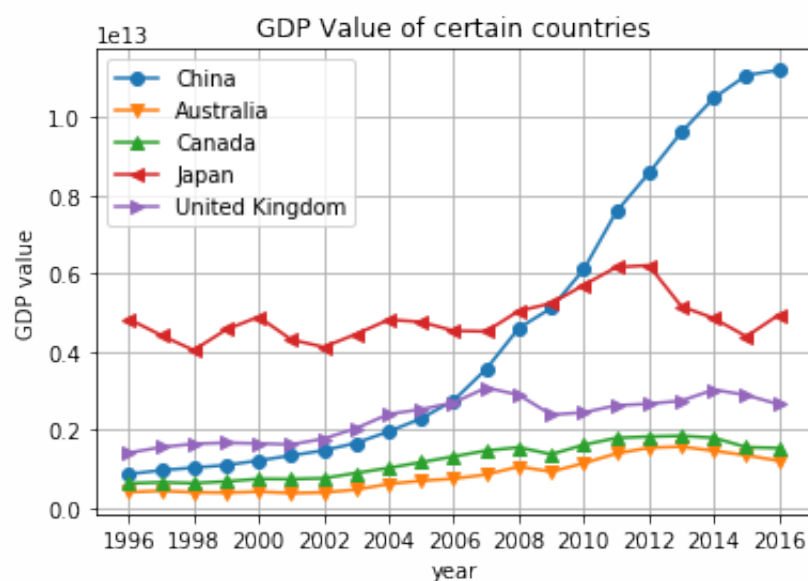
题目三：GDP数据可视化

题目描述

Visualize the GDP values and changes of GDP of the countries for the past 20 years. The data can be download from <http://www.gapminder.org>.

解答

```
1 import pandas as pd
2 import numpy as np
3 from matplotlib import pyplot as plt
4
5 ##### 数据生成 #####
6 GDP_raw = pd.read_csv("GDP.csv", index_col=0, skiprows=4)
7
8 year = list(range(1996,2017))
9 cty = ["China", "Australia", "Canada", "Japan", "United Kingdom"]
10 GDP = GDP_raw[[str(x) for x in year]].loc[cty]
11
12 ##### 作图过程 #####
13 fig = plt.figure()
14 markerSymbols = ['o-', 'v-', '^-', '<-', '>-', 's-', 'p-', 'P-']
15
16 for i in range(len(cty)):
17     plt.plot(year, GDP.values[i], markerSymbols[i])
18
19 plt.title("GDP Value of certain countries")
20 plt.ylabel("GDP value")
21 plt.xlabel("year")
22 plt.xticks(range(1996,2017,2))
23 plt.legend(cty, loc='upper left')
24
25 plt.grid()
26 plt.show()
```



题目四：直方图函数

题目描述

Program a function for computing the histogram of a data set of N dimension. The function should include the definition/setting of (1) number of bins, (2) width(s) of bin(s), (3) frequencies of each bin, and (4) the dimension of the input data. Test the function using a dataset and visualize the result using bar plot. Submit your Python code and test data.

解答

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4
5 def my_histogram(data, binNum, dim=None):
6     # 若为一维或高维 (>1) 则选取一维做频率直方图
7     data = data if dim is None else data[dim]
8     dataMin, dataMax = min(data), max(data)
9
10    # 构造直方图位置列表
11    bins = np.linspace(dataMin, dataMax, binNum)
12    if binNum == 1:
13        N = np.array([len(data)])
14        width = dataMax - dataMin
15    else:
16        N = np.zeros(binNum, dtype=np.int64)
17        width = (dataMax - dataMin) / (binNum - 1)
18        for d in data:
19            # 对小范围内的数据统计频率
20            N[int((d - dataMin)//width)] += 1
21
22        # 把最大值归为最后一个区间
23        N[-2] += 1
24    return (N[:-1], bins, width)
25
26
27 ##### 数据生成 #####
28 # 让数据采样足够大以此展现出数据分布的随机（平均）性
29 n = 10000
30 data = np.array([np.random.rand(n),
31                  np.random.rand(n),
32                  np.random.rand(n)])
33
34 ##### 作图过程 #####
```

```

35 # 选取第2维
36 N, pos, width = my_histogram(data, 10, dim=2)
37
38 # 设置图片大小 方便横坐标足够显示清楚
39 plt.figure(figsize=(12,7))
40
41 # 设置横坐标为数据区间
42 xtick = ["%.2f" % pos[i] + ' ~ ' + "%.2f" % pos[i+1] for i in
43          range(len(pos)-1)]
44
45 plt.bar(xtick, N, width=width*5, color='b', alpha=0.5)
46
47 plt.title("Histogram plot")
48 plt.show()

```

