



用 Python 做统计分析

作者：陈震

组织：西南大学

时间：September 14, 2019

版本：1.00



Victory won't come to us unless we go to it. — M. Moore

目 录

1	Python 概述与安装	1
1.1	Python 的发展历史与应用	1
1.2	安装 Python	2
1.2.1	Python + Pycharm	2
1.2.2	Anaconda + Spyder	4
1.3	统计学常用 Python 包的安装	4
2	Python 的操作基础	5
2.1	Python 的常用数据类型	5
2.2	Python 的常用数据结构	5
2.3	Python 中定义函数	5
2.4	Python 读取数据	5
3	Python 中的 Numpy 与 Pandas	6
3.1	Numpy 的常用操作	6
3.2	Pandas 的常用操作	6
4	统计数据的图形描述	7
4.1	plot 函数的基本用法	7
4.2	常用图形的绘制	10
4.2.1	线图	10
4.2.2	散点图	11
4.2.3	条形图	16
4.2.4	直方图	20
4.2.5	饼图	22
4.2.6	箱线图	25
4.3	用 pandas 调用 plot 函数画图	25
4.4	画子图函数 subplot	25
5	数据的概括性度量	26
5.1	均值、方差	26
5.2	众数、中位数、分位数	26
5.3	偏态、峰态	26
6	参数估计与假设检验	27

7 方差分析	28
8 线性回归	29
9 非线性回归	30
10 主成分分析	31
11 因子分析	32
12 时间序列分析和预测	33

第 1 章 Python 概述与安装

1.1 Python 的发展历史与应用

Python 的诞生时间非常早，第一版 Python 发行于 1991 年，甚至比 Java 的历史都早。在大部分时间内，Python 一直作为一个小众的编程语言，并没有大规模流行起来。近些年来，随着大数据技术、人工智能技术的发展和普及，“简洁、具有良好扩展性”的 Python 非常契合大数据与人工智能技术对编程语言的要求，超越其他编程语言迅速崛起。Python 分别在 2010 年，2018 年被全球知名的编程语言流行度排行榜网站 TIOBE 评为“年度最佳编程语言”，并长期位于各年或各月流行编程语言排行榜的前三名。中国的网络上甚至产生了一句流行语：“人生苦短，我用 Python”。

Python 的创始人为荷兰人吉多·范罗苏姆（Guido van Rossum），在 1989 年的圣诞节期间，吉多·范罗苏姆为了打发时间，决心开发一个新的脚本解释程序，作为 ABC 语言的继承，于是 Python 诞生了。之所以选择 Python 作为名字，是由于吉多·范罗苏姆非常喜欢一部 BBC 电视剧—Monty Python’s Flying Circus（中文译名：蒙提·派森的飞行马戏团）。Python 的英文意思为蟒蛇，这也是为什么一些 Python 相关的软件或书籍用蟒蛇作为图标的原因。

Python 的设计哲学是“优雅”、“明确”、“简单”，Python 可以说是语法功能最简单的编程语言之一。Python 具有良好的可扩展性，Python 提供了丰富的接口和工具，以便程序员能够轻松地使用 C、C++、Cython 来编写扩展模块，因此，有很多人把 Python 作为一种“胶水语言”使用。Python 是完全面向对象的语言，函数、模块、数字、字符串都是对象，并且完全支持继承、重载、派生、多重继承等，Python 还支持重载运算符。并且随着计算机速度运算速度越来越快，Python 运算速度慢的缺点逐渐可以忽略，因此，许多网站开始使用 Python 开发，许多专业软件也开始支持 Python 调用。

Python 的应用范围包括 Web 开发网站与网络爬虫，GUI 开发桌面软件，科学计算等。使用 Python 编写的著名应用包括：

- Youtube — 视频分享网站
- Dropbox — 文件分享软件
- Instagram — 图片分享软件
- 豆瓣网 — 国内著名的图书、唱片、电影评论网站
- 知乎网 — 国内著名的问答网站
- 果壳网 — 国内著名的泛科技主题网站

目前的主要 Python 版本为 Python3，本书的所有 Python 程序都是 Python3 版本。

1.2 安装 Python

使用 Python 时一般有两种方式：一是下载 Python 和编辑器 Pycharm，二是下载 Anaconda（自带 Python）和编辑器 Spyder。从实践经验来看，Pycharm 功能比 Spyder 更多，更加人性化，笔者优先推荐第一种方式。

1.2.1 Python + Pycharm

1.2.1.1 下载 python

Python 的官方网址为：<https://www.python.org>。打开网址后，从 Download 里面可以找到不同操作系统的 python 版本下载。下面我们以 Windows 系统为例讲述 Python 的安装。



下载完成后，双击 exe 文件安装，切记要勾选 Add Python 3.7 to PATH，然后点击 Install Now（默认安装目录，默认安装 pip 等）或 Customize installation（可以自己设置安装目录等）都可以。



为了测试 Python 是否安装成功，打开命令行窗口（win10 系统中敲 win 键，再输入

cmd, 然后回车), 输入 python, 若能显示 Python 的版本信息, 则代表 python 已正确安装。

```
命令提示符 - python
Microsoft Windows [版本 10.0.17763.737]
(c) 2018 Microsoft Corporation。保留所有权利。

C:\Users\chen>python
Python 3.7.3 (v3.7.3:ef4ec6ed12, Mar 25 2019, 22:22:05) [MSC v.1916 64 b
type "help", "copyright", "credits" or "license" for more information.
>>>
```

1.2.1.2 下载 Pycharm

Pycharm 的官方网址: <http://www.jetbrains.com/pycharm>, 点击 DOWNLOAD 开始安装。

Download PyCharm

Windows

macOS

Linux

Professional

For both Scientific and Web Python development. With HTML, JS, and SQL support.

DOWNLOAD

Free trial

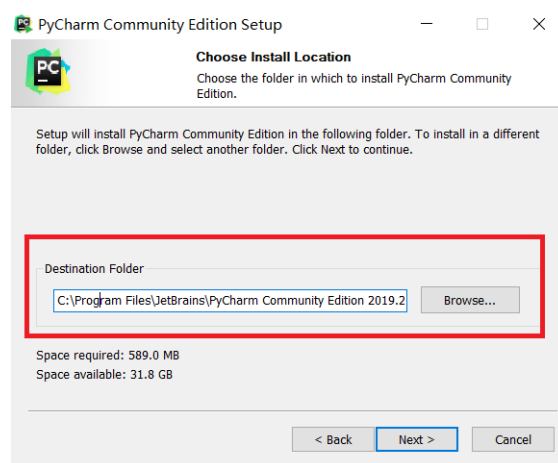
Community

For pure Python development

DOWNLOAD

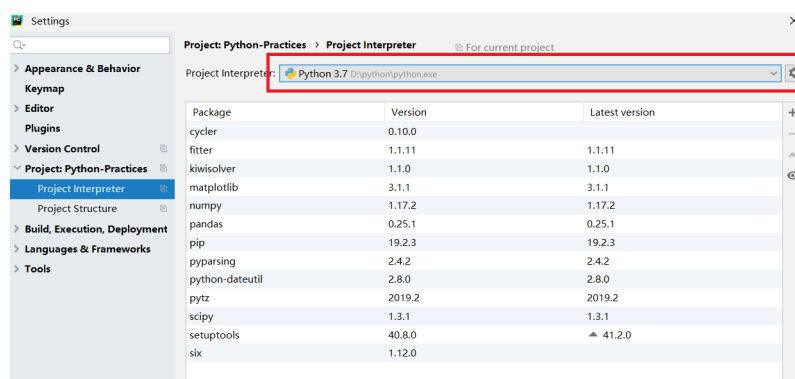
Free, open-source

一般选择免费版本下载, 下载完成后双击 exe 进行安装。



可以通过图中的 Browse 更改安装目录, 其他步骤按照默认设置即完成安装。

完成安装后，设置 Pycharm 中的 Python 解释器为我们安装的 Python 版本。打开 Pycharm，依次点击 file—settings—project interpreter，将 interpreter 设置为 Python 安装路径中的 python.exe（笔者电脑的位置为 D:\python:\python.exe）。



1.2.2 Anaconda + Spyder

1.3 统计学常用 Python 包的安装

第 2 章 Python 的操作基础



2.1 Python 的常用数据类型

2.2 Python 的常用数据结构

2.3 Python 中定义函数

2.4 Python 读取数据

第 3 章 Python 中的 Numpy 与 Pandas



3.1 Numpy 的常用操作

3.2 Pandas 的常用操作

第 4 章 统计数据的图形描述

使用图形描述统计结果是应用统计的基本技能之一。人类在接受信息时，大脑皮层总是优先提取可视化的图形，其次才是具体的文字内容。一些漂亮的图形会使得读者或观众迅速得知我们想要表达的信息。因此不论是写论文，还是汇报展示，做到图文并茂非常重要。

由于 Python 具有良好的可扩展性，它在画图方面具有大量的绘图包，包括模仿 Matlab 绘图风格的 `matplotlib`，画地理图形的 `geoplotlib`，能在浏览器中生成优美图形的 `Pyechart` 等。本章将展示如何用 Python 的 `matplotlib` 宏包来画出各种各样的统计图形。

4.1 `plot` 函数的基本用法

使用 `matplotlib` 包画图时，我们一般加载里面的 `pyplot`，并命名为 `plt`，然后使用 `plot` 函数画图。

```
# 导入 matplotlib 中的 plot，并命名为常用名 plt
import matplotlib.pyplot as plt
```

例如，下面的代码画出正弦函数 $y = \sin(x)$ 的图形。

```
# 导入宏包
import matplotlib.pyplot as plt
import numpy as np

# 生成数据
x = np.arange(0, 10, 0.1) # 横坐标数据为从0到10之间，步长为0.1的等差数组
y = np.sin(x) # 纵坐标数据为 x 对应的 sin(x) 值

# 生成图形
plt.plot(x, y)

# 显示图形
plt.show()
```

生成的图形如图 4.1 所示。利用 `plot` 函数，我们可以对图形进行更多精细的设置，官方的详细文档可以参看：https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.plot.html。Plot 函数的基本语法是：

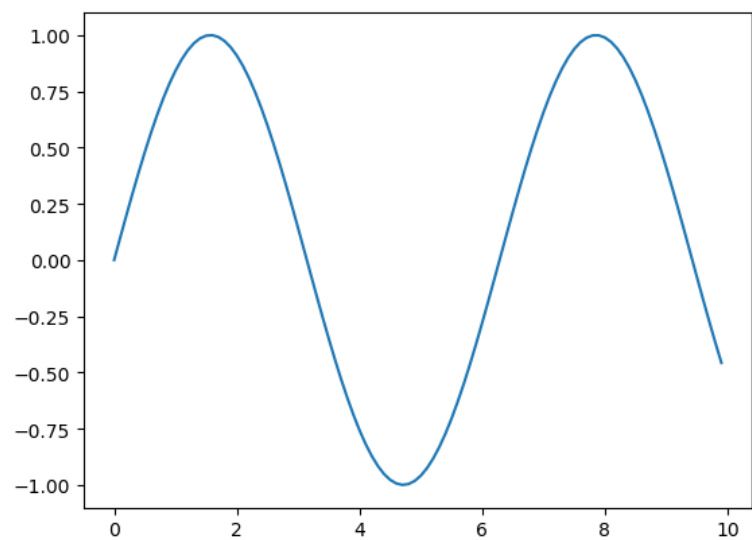


图 4.1: $y = \sin(x)$ 的 Python 图像

plot 函数的语法

plot([x], y, [fmt], **kwargs)

[x]

可选参数，横坐标轴数据

y

纵坐标轴数据

[fmt]

可选参数，字符串，定义图形的基本样式：颜色，点形，线形

****kwargs**

不定长的关键字参数，用字典形式设置图形的其他属性

[fmt] 的常用代码（包括颜色代码、点形代码、线形代码），由表 4.1 所示。

表 4.1: Plot 函数中的颜色、点形、线性代码

(a) Plot 函数的颜色代码		(b) Plot 函数的点形代码		(c) Plot 函数的线形代码	
颜色代码	颜色	线形代码	线形	线形代码	线形
'b'	蓝色	'o'	实心圆形	'-'	实线
'r'	红色	'.'	点形	'--'	虚线
'g'	绿色	'+'	十字形	'-.'	折线
'k'	黑色	'*'	星号	'.'	点线
'w'	白色	'x'	加号		
'y'	黄色		叉号		

****kwargs** 的常用设置包括线条的粗细 **linewidth**（或写成 **lw**），图像标签 **label** 等。下面一些 **plot** 函数的代码展示了 **[x]**，**[fmt]**，****Kwargs** 的一些可选用法。



```
>>> plot(x, y)      # 根据横坐标数据 x 与纵坐标数据 y 画图, 采用默认的颜色、点形与
                    线性
>>> plot(y)        # 据纵坐标数据 y 画图, 横坐标数据默认为从 0 到 N-1, 步长为 1
                    的等差数组
>>> plot(x, y, 'bo') # 颜色为蓝色('b')、点形为圆('o')
>>> plot(y, 'g-.')  # 颜色为绿色('g'), 线形为折线('-.')
>>> plt.plot(x, y, 'yo:', label='y=sin(x)', lw=2) # 颜色为黄色('y'), 点形为圆形
                    ('o'), 线形为虚线(':'), label 内容为 'y=sin(x)', 线条宽度为 2
```

如果我们想自定义坐标轴的标题, 坐标轴的刻度, 坐标轴刻度的范围, 设置图形标题, 添加图例时, 可以通过设置 `pyplot` 函数中的 `xlabel` (横坐标轴标题), `ylabel` (纵坐标轴标题), `xticks` (横坐标轴刻度), `yticks` (纵坐标轴刻度), `title` (图形标题), `grid` (显示网格), `legend` (显示图例) 等属性来实现。经过自定义设置, 对图 4.1 的代码进行一下修改, 生成图 4.2。

```
# 导入宏包
import matplotlib.pyplot as plt
import numpy as np

# 这两行代码使得 pyplot 画出的图形中可以显示中文
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 生成数据
x = np.arange(0, 10, 0.5)
y = np.sin(x)

# 生成图形
plt.plot(x, y, 'go:', label='y=sin(x)', linewidth=2) # 颜色绿色, 点形圆形, 线性
                                                       虚线, 设置图例显示内容, 线条宽度为2

plt.ylabel('y') # 横坐标轴的标题
plt.xlabel('x') # 纵坐标轴的标题
plt.xticks(np.arange(0, 11, 1)) # 设置横坐标轴的刻度为 0 到 10 的数组
plt.ylim([-2, 2]) # 设置纵坐标轴范围为 -2 到 2
plt.legend() # 显示图例, 图例中内容由 label 定义
plt.grid() # 显示网格
plt.title('我的第一个 Python 图形') # 图形的标题

# 显示图形
plt.show()
```

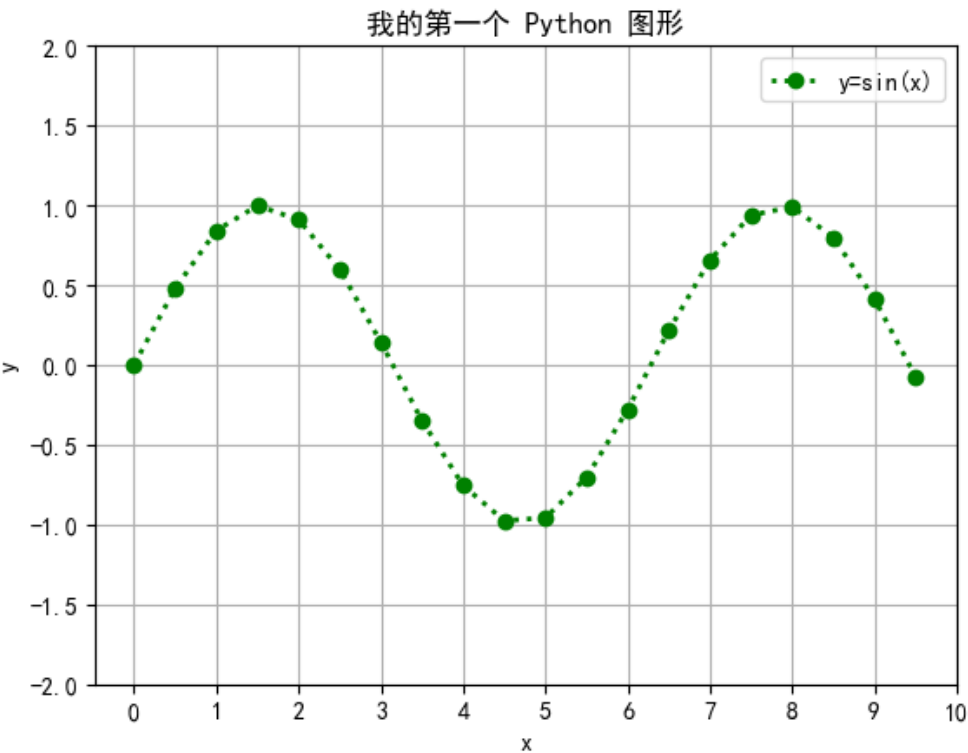


图 4.2: 自定义坐标轴标题、刻度、范围，显示图例、网格与图形标题

4.2 常用图形的绘制

4.2.1 线图

在统计学中，线图一般用来表示时间序列数据，线图称为趋势图（run chart），因为从线图中可以清晰地看出数据虽时间变化的趋势。表 4.2 是我国近 10 年的 GDP 增长率，以及三大产业在近 10 年的增长率。

表 4.2: 我国近 10 年的经济增长率，数据来源：国家统计局

时间	GDP 增长率	第一产业增长率	第二产业增长率	第三产业增长率
2009 年	9.40	4.00	10.30	9.60
2010 年	10.60	4.30	12.70	9.70
2011 年	9.60	4.20	10.70	9.50
2012 年	7.90	4.50	8.40	8.00
2013 年	7.80	3.80	8.00	8.30
2014 年	7.30	4.10	7.40	7.80
2015 年	6.90	3.90	6.20	8.20
2016 年	6.70	3.30	6.30	7.70
2017 年	6.80	4.00	5.90	7.90
2018 年	6.60	3.50	5.80	7.60

在画图时，横坐标轴数据为年份，纵坐标轴数据分别为 GDP 增长率，第一产业增长



率，第二产业增长率，第三产业增长率。为了将四个纵坐标轴数据显示在一个图形上，可以用四个 plot 函数进行划线。Python 画图的代码为：

```
import matplotlib.pyplot as plt

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 输入纵坐标轴数据与横坐标轴数据
gdp_rate = [9.4, 10.6, 9.6, 7.9, 7.8, 7.3, 6.9, 6.7, 6.8, 6.6]
first_industry_rate = [4.0, 4.3, 4.2, 4.5, 3.8, 4.1, 3.9, 3.3, 4.0, 3.5]
second_industry_rate = [10.3, 12.7, 10.7, 8.4, 8.0, 7.4, 6.2, 6.3, 5.9, 5.8]
third_industry_rate = [9.6, 9.7, 9.5, 8.0, 8.3, 7.8, 8.2, 7.7, 7.9, 7.6]
years = [2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018]

# 4 个 plot 函数画出 4 条线，线形为折线，每条线对应各自的标签 label
plt.plot(years, gdp_rate, '-.', label='GDP增长率')
plt.plot(years, first_industry_rate, '-.', label='第一产业增长率')
plt.plot(years, second_industry_rate, '-.', label='第二产业增长率')
plt.plot(years, third_industry_rate, '-.', label='第三产业增长率')

plt.xticks(years) # 设置横坐标刻度为给定的年份
plt.xlabel('年份') # 设置横坐标轴标题
plt.legend() # 显示图例，即每条线对应 label 中的内容
plt.show() # 显示图形
```

显示效果如图 4.3 所示。

4.2.2 散点图

散点图是数据点在直角坐标系平面上的分布图，在统计学的回归分析与预测中经常用到。用横轴代表变量 x ，纵轴代表变量 y ，每组数据 (x_i, y_i) 在坐标系中用一个点表示。做散点图要用到 pyplot 中的 scatter 函数，该函数的基本语法为¹：

¹在 matplotlib 宏包中，0 到 1 之间的浮点数产生相应的 RGB 颜色

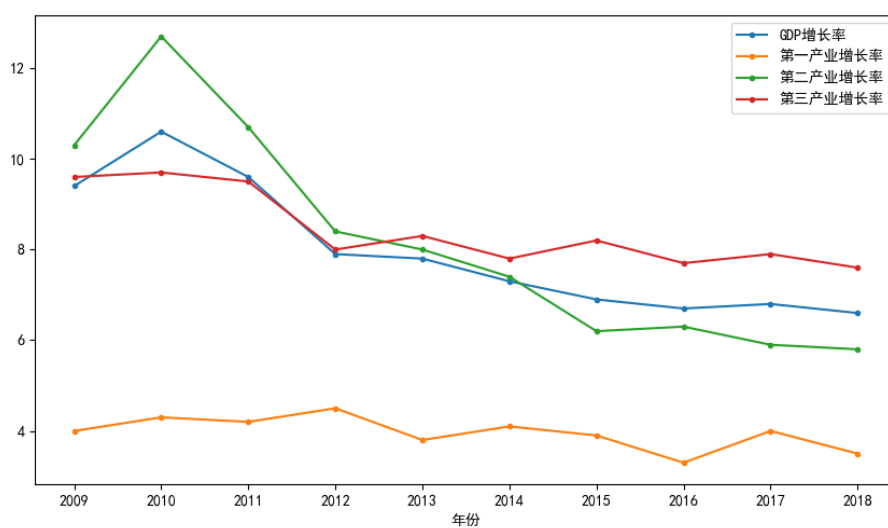


图 4.3: 近十年我国的 GDP 以及各产业增长率

scatter 函数的语法

scatter(x, y, [s], [c], **kwargs)

- x** 横坐标轴数据
- y** 纵坐标轴数据
- [s]** 可选参数，一个数或一个数组，设置每个散点的大小
- [c]** 可选参数，一个数或一个数组，设置每个散点的颜色
- **kwargs** 不定长的关键字参数，用字典形式设置图形的其他属性

****kwargs** 中常设置的是不透明度属性 **alpha**，其大小为 0 到 1 之间的浮点数。假设某个农产品的产量与温度和降雨量的关系如表 4.3 所示。

表 4.3: 某农产品的产量与温度和降雨量的关系

产量	温度	降雨量
1125	6	25
1725	8	40
2250	10	58
2875	13	68
2900	14	110
3750	16	98
4125	21	120

作出产量与温度的散点图的 Python 代码与图形如下：

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



```
import numpy as np

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 输入产量与温度数据
production = [1125, 1725, 2250, 2875, 2900, 3750, 4125]
tem = [6, 8, 10, 13, 14, 16, 21]

colors = np.random.rand(len(tem)) # 颜色数组
plt.scatter(tem, production, s=200, c=colors) # 画散点图, 大小为 200
plt.xlabel('温度') # 横坐标轴标题
plt.ylabel('产量') # 纵坐标轴标题
plt.show()
```

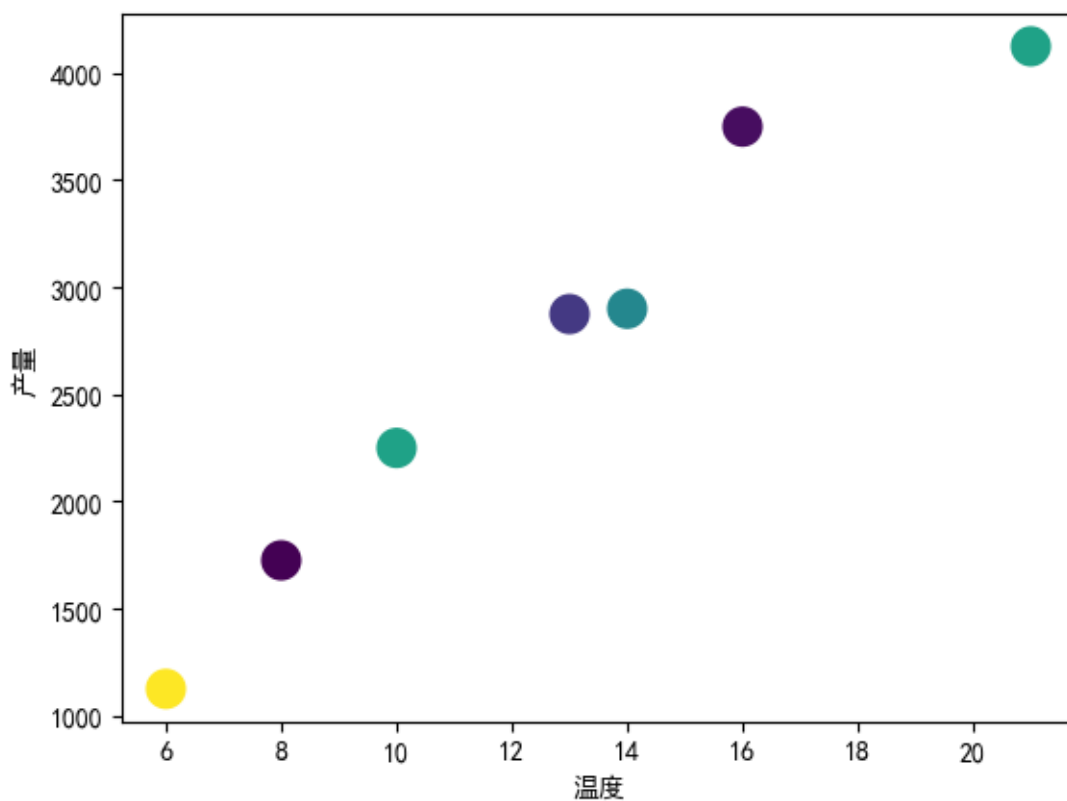


图 4.4: 产量与温度的散点图

若将散点大小的数据换为第三个变量的数值, 则可以作出反映三个变量关系的气泡图。下面的代码和图形做出了一个气泡图。图 4.5 反映了产量与温度、降雨量的关系: 温度数值在横坐标轴, 降雨量数值在纵坐标轴, 降雨量的大小用气泡的大小表示。

```
import matplotlib.pyplot as plt
import numpy as np

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 输入产量与温度数据
production = [1125, 1725, 2250, 2875, 2900, 3750, 4125]
tem = [6, 8, 10, 13, 14, 16, 21]
rain = [25, 40, 58, 68, 110, 98, 120]

colors = np.random.rand(len(tem)) # 颜色数组
size = production
plt.scatter(tem, rain, s=size, c=colors, alpha=0.6) # 画散点图, alpha=0.6 表示不
    透明度为 0.6
plt.ylim([0, 150]) # 纵坐标轴范围
plt.xlim([0, 30]) # 横坐标轴范围
plt.xlabel('温度') # 横坐标轴标题
plt.ylabel('降雨量') # 纵坐标轴标题
plt.show()
```

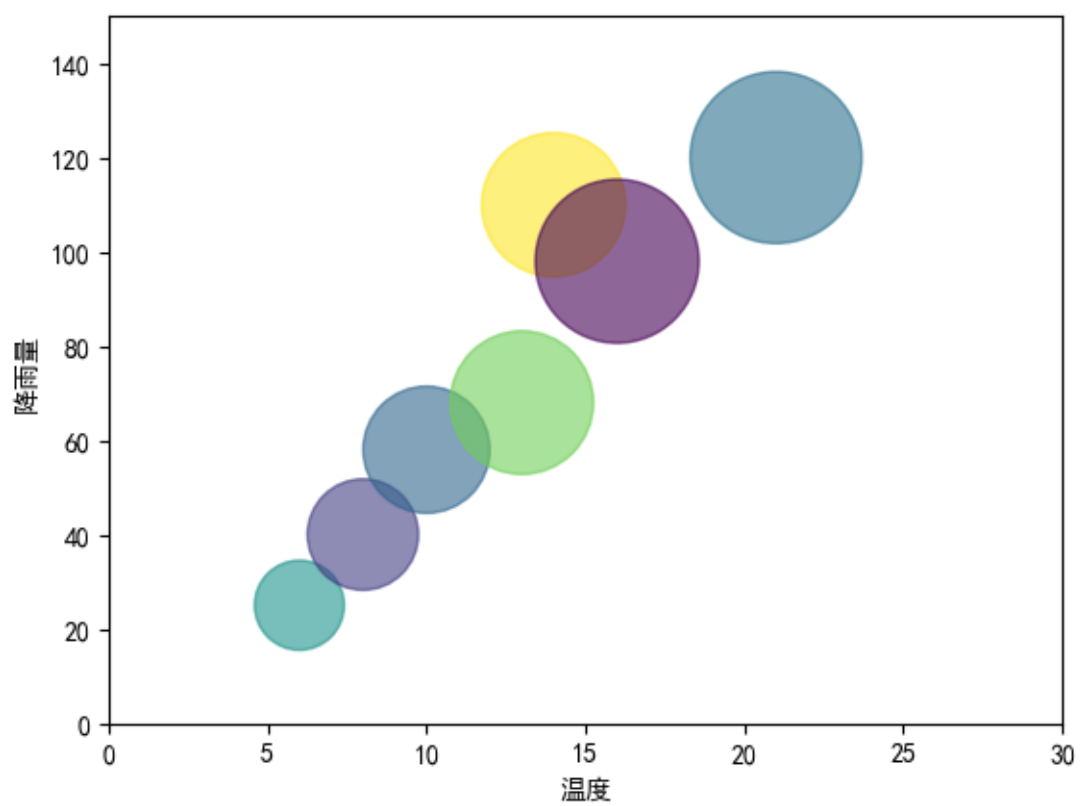


图 4.5: 产量与温度、降雨量关系的气泡图

4.2.3 条形图

条形图 (bar chart), 也称为柱状图, 是一种以长方形的长度为变量的统计图表, 长方形的长度与它所对应的变量数值呈一定比例。画条形图要用到 `pyplot` 中的 `bar` 函数, 该函数的基本语法为:

bar 函数的语法

bar(x, height, [width], **kwargs)

x 数组, 每个条形的横坐标
height 一个数或一个数组, 条形的高度
[width] 可选参数, 一个数或一个数组, 条形的宽度, 默认为 0.8
****kwargs** 不定长的关键字参数, 用字典形式设置条形图的其他属性

`**kwargs` 中常设置的参数包括图形标签 `label`, 颜色标签 `color`, 不透明度 `alpha` 等。假设某项针对男女大学生购买饮用水爱好的调查结果如下表:

表 4.4: 男女大学生购买饮用水爱好调查表

买水选择	男	女
碳酸饮料	6	9
绿茶	7	4
矿泉水	6	4
果汁	1	5
其他	2	6
总计	22	28

画出男生饮用水情况的直方图, 代码如下:

```
import matplotlib.pyplot as plt

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

waters = ('碳酸饮料', '绿茶', '矿泉水', '果汁', '其他')
buy_number = [6, 7, 6, 1, 2]

plt.bar(waters, buy_number)
plt.title('男性购买饮用水情况的调查结果')

plt.show()
```

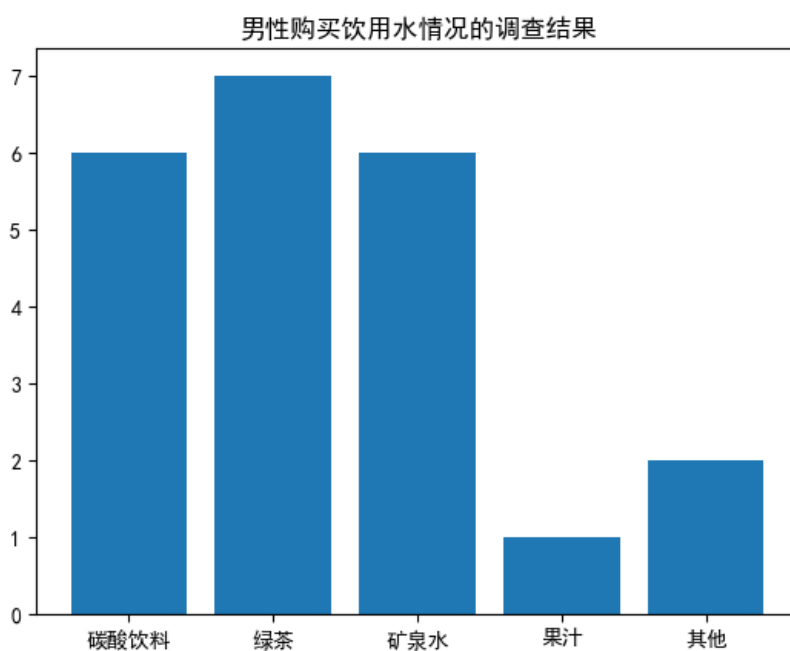


图 4.6: 男性购买饮用水条形图

显示的图形:

若要生成横的条形图, 则可以使用 `barh` 函数, 其语法与 `bar` 函数非常类似。

barh 函数的语法

bar(x, width, [height], **kwargs)

- y** 数组, 每个条形的纵坐标
- width** 一个数或一个数组, 条形的宽度
- [height]** 可选参数, 一个数或一个数组, 条形的高度, 默认为 0.8
- **kwargs** 不定长的关键字参数, 用字典形式设置条形图的其他属性

对应的代码与图形为:

```
import matplotlib.pyplot as plt

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

waters = ('碳酸饮料', '绿茶', '矿泉水', '果汁', '其他')
buy_number = [6, 7, 6, 1, 2]

plt.barh(waters, buy_number) # 横放条形图函数 barh
```

```
plt.title('男性购买饮用水情况的调查结果')

plt.show()
```

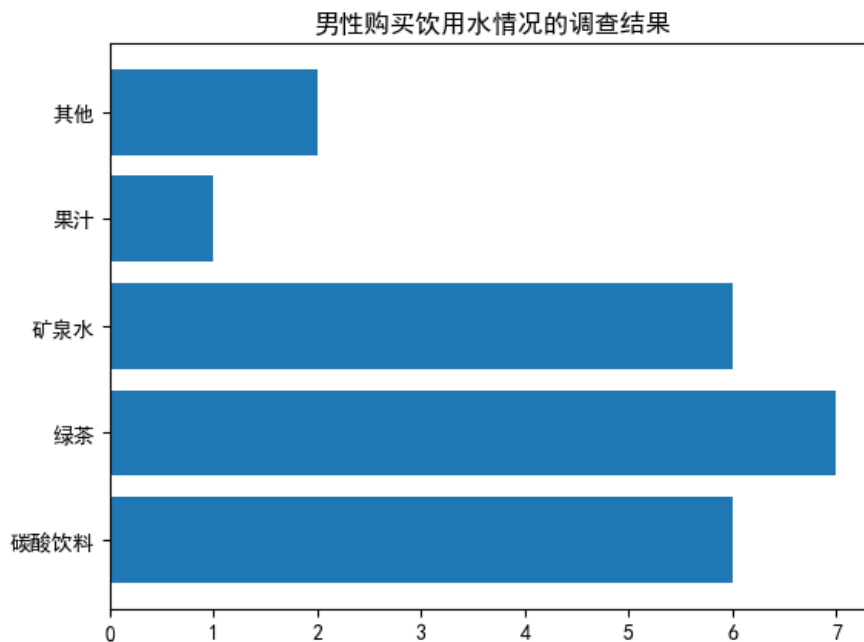


图 4.7: 男性购买饮用水条形图（条形横放）

若要将男生与女生的调查情况画出两个条形图一块显示，则可以使用 `bar` 或 `barh` 函数两次，并调整 `bar` 或 `barh` 函数的条形图位置坐标以及相应刻度，使得两组条形图能够并排显示。

```
import matplotlib.pyplot as plt
import numpy as np

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 输入统计数据
waters = ('碳酸饮料', '绿茶', '矿泉水', '果汁', '其他')
buy_number_male = [6, 7, 6, 1, 2]
buy_number_female = [9, 4, 4, 5, 6]

bar_width = 0.3 # 条形宽度
index_male = np.arange(len(waters)) # 男生条形图的横坐标
index_female = index_male + bar_width # 女生条形图的横坐标
```

```
# 使用两次 bar 函数画出两组条形图
plt.bar(index_male, height=buy_number_male, width=bar_width, color='b', label='
    男性')
plt.bar(index_female, height=buy_number_female, width=bar_width, color='g',
    label='女性')

plt.legend() # 显示图例
plt.xticks(index_male + bar_width/2, waters) # 让横坐标轴刻度显示 waters 里的饮
    用水, index_male + bar_width/2 为横坐标轴刻度的位置
plt.ylabel('购买量') # 纵坐标轴标题
plt.title('购买饮用水情况的调查结果') # 图形标题

plt.show()
```

生成图 4.8。

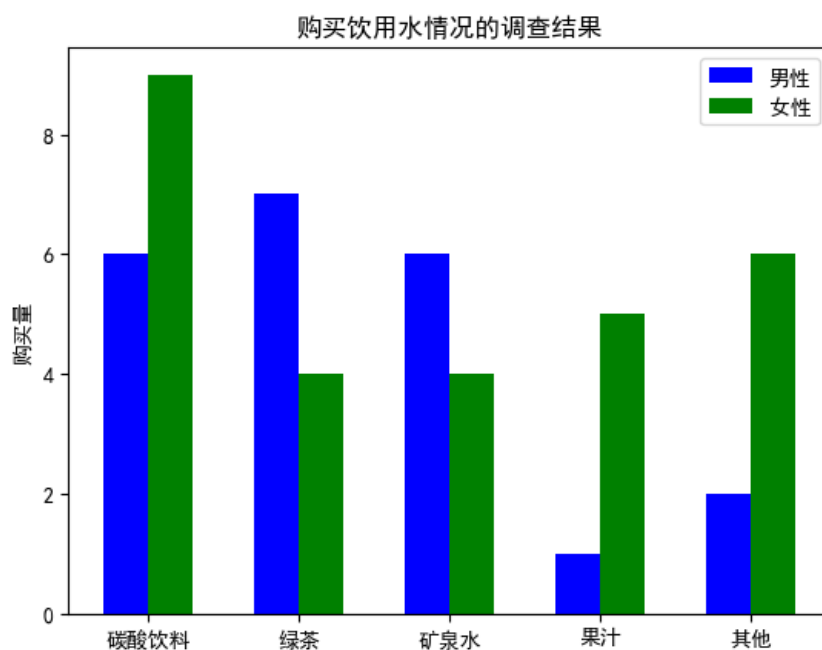


图 4.8: 男女购买饮用水条形图

4.2.4 直方图

直方图（histogram）虽然在样式上类似条形图，但它们的作用不一样。直方图用不同的矩形表示频数，常用来观察一组数据的概率分布。在直角坐标中，用横轴表示数据分组，纵轴表示频数或频率，各组与相应的频数就形成了一个矩形，即直方图。

画直方图用到 pyplot 中的 hist 函数，它的基本语法为：

hist 函数的语法

[n, bins, patches] = hist(x, [bins], **kwargs)

输入值：

x 数组，需要绘制直方图的数值
[bins] 可选参数，数据的组数，若不指定则 hist 函数默认计算一个组数
****kwargs** 不定长的关键字参数，用字典形式设置条形图的其他属性

返回值：

n 一个数组，每组直方图的频数（或概率密度）
bins 一个数组，直方图的边界数值
patches 一个对象组，每个对象代表直方图的矩形

****kwargs** 中常用来设置的属性包括直方的边界线颜色 **edgecolor**，不透明度 **alpha** 等。需要注意的是，hist 函数的三个返回值一般用来继续对数据进行分析，并不是必须要写或必须要使用的。下面举例说明该函数，例如某个饭店每天接待的顾客数有以下记录值：

表 4.5: 某饭店每天接待的顾客数的记录值

141	159	166	172	177	182	188	196	203	214
143	160	167	173	177	183	189	196	203	215
144	160	168	173	178	184	189	196	205	218
149	161	168	174	178	185	189	196	206	223
150	161	168	174	178	186	190	196	207	225
152	162	170	174	179	186	190	197	208	226
153	163	171	175	179	187	191	197	209	228
153	163	171	175	179	187	192	198	210	233
154	164	172	175	180	187	194	198	210	233
155	165	172	175	180	187	194	200	211	234
156	165	172	176	181	188	195	201	211	234
158	165	172	176	182	188	195	202	213	237

在 Phthon 中编写代码，并画出对应的直方图。

```
import matplotlib.pyplot as plt

x = [141, 159, 166, 172, 177, 182, 188, 196, 203, 214,
     143, 160, 167, 173, 177, 183, 189, 196, 203, 215,
     144, 160, 168, 173, 178, 184, 189, 196, 205, 218,
     149, 161, 168, 174, 178, 185, 189, 196, 206, 223,
     150, 161, 168, 174, 178, 186, 190, 196, 207, 225,
     152, 162, 170, 174, 179, 186, 190, 197, 208, 226,
     153, 163, 171, 175, 179, 187, 191, 197, 209, 228,
     153, 163, 171, 175, 179, 187, 192, 198, 210, 233,
     154, 164, 172, 175, 180, 187, 194, 198, 210, 233,
     155, 165, 172, 175, 180, 187, 194, 200, 211, 234,
     156, 165, 172, 176, 181, 188, 195, 201, 211, 234,
     158, 165, 172, 176, 182, 188, 195, 202, 213, 237]

plt.hist(x, edgecolor='k', alpha=0.35) # 设置直方边线颜色为黑色，不透明度为 0.35
plt.show()
```

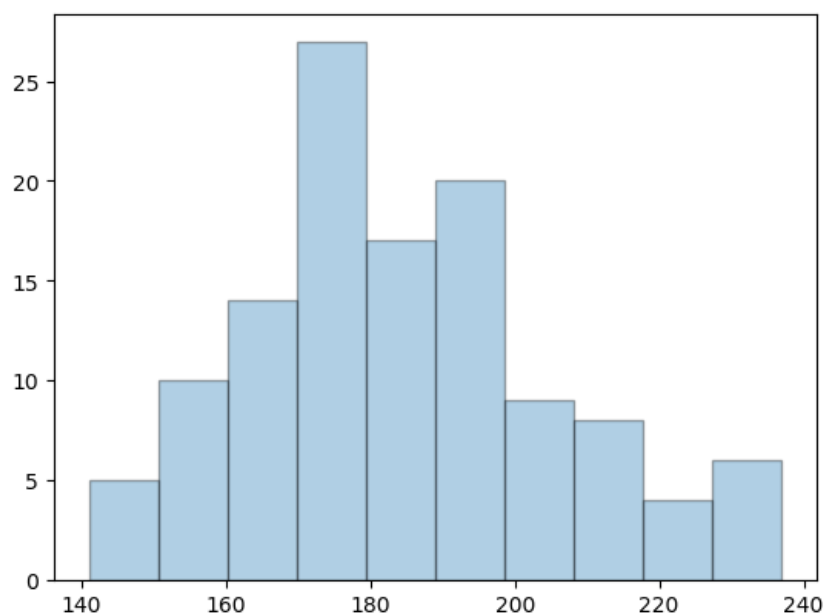


图 4.9: 每天接待顾客数分布的直方图

从图 4.9 可以看出，每天的顾客数大致符合一个正态分布。

4.2.5 饼图

饼图（pie char）是一个划分为几个扇形的圆形统计图表，一般用于描述频率或百分比之间的相对关系。在饼图中，每个扇区的弧长（以及圆心角和面积）的大小与其所表示的数量呈固定比例。

画饼图一般使用 pyplot 中的 pie 函数，它的基本语法如下：

pie 函数的语法

pie(x, [expode], [labels], [autopct], **kwargs)

x	数组，每个扇区的比例
[expode]	可选参数，数组，每个扇区突出的大小
[labels]	可选参数，字符串数组，每个扇区的标签
[autopct]	可选参数，字符串或函数，每个扇区显示的数字样式
**kwargs	不定长的关键字参数，用字典形式设置条形图的其他属性

****kwargs** 中常用来设置的属性包括是否显示扇区的阴影 shadow，扇区的起始角度 startangle 等。假设有下列调查数据：

表 4.6: 购水类型的调查数据

购水类型	人数
碳酸饮料	15
绿茶	11
矿泉水	10
其他	8
果汁	6

饼图的代码与图形：

```
import matplotlib.pyplot as plt
import numpy as np

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

labels = ['果汁', '矿泉水', '绿茶', '其他', '碳酸饮料']
x = [6, 10, 11, 8, 15]
explode = [0, 0.1, 0, 0, 0] # 突出显示第二个扇区

plt.pie(x, explode=explode, labels=labels, autopct='%2f%%', shadow=True,
        startangle=90)
```

```
plt.legend() # 显示标签
plt.axis('equal') # 让图形和坐标轴相等，这样避免饼图为椭圆形
plt.show()
```

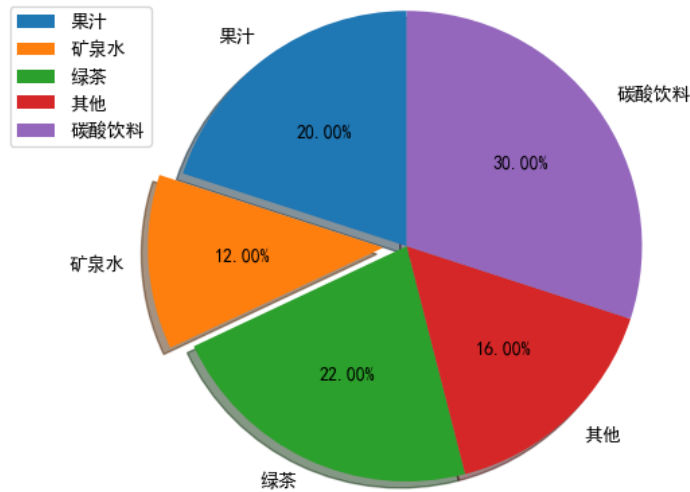


图 4.10: 调查结果的饼图

若在饼图上既要显示比率，又要显示实际调查数字，则可以通过自定义 `autopct` 函数来实现：

```
import matplotlib.pyplot as plt
import numpy as np

# 这两行代码解决 plt 中文显示的问题
plt.rcParams['font.sans-serif'] = ['SimHei']
plt.rcParams['axes.unicode_minus'] = False

# 自定义 autopct 函数
def func(pct, allvals):
    absolute = int(round(pct*np.sum(allvals)/100.0))
    return "{:.1f}%\n({:d})".format(pct, absolute)

labels = ['果汁', '矿泉水', '绿茶', '其他', '碳酸饮料']
x = [6, 10, 11, 8, 15]
explode = [0, 0.1, 0, 0, 0] # 突出显示第二个扇区

plt.pie(x, explode=explode, labels=labels, autopct=lambda pct: func(pct, x), #
```

```
利用 lambda 定义 pct 函数
    shadow=True, startangle=90)
plt.legend() # 显示标签
plt.axis('equal') # 让图形和坐标轴相等，这样饼图会更好看
plt.show()
```

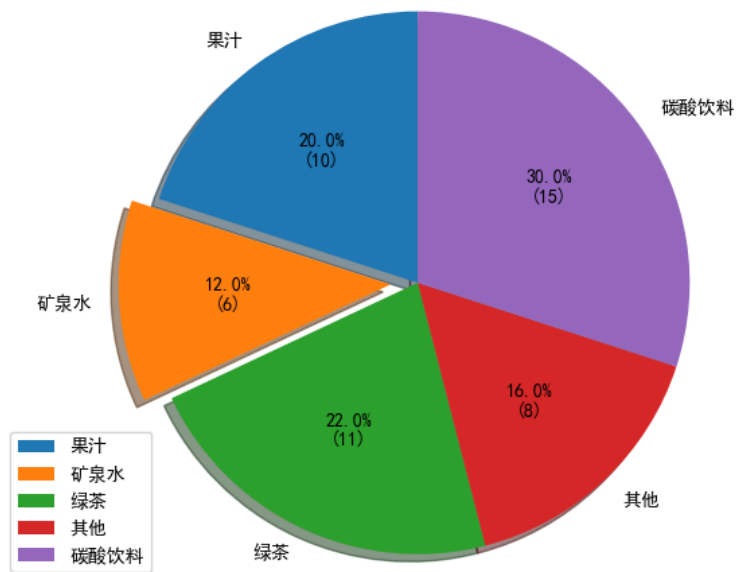


图 4.11: 调查结果的饼图（同时显示比例与原始数据）

4.2.6 箱线图

箱线图（box plot）在学术论文中经常出现，是一种用来显示一组数据分散情况资料的统计图，由数据的最大值、最小值、中位数、两个四分位数，共五个特征绘制而成。

表 4.7: 学生成绩统计表

	学生编号										
课程名称	1	2	3	4	5	6	7	8	9	10	11
英语	76	90	97	71	70	93	86	83	78	85	81
西方经济学	93	81	76	88	66	79	83	92	78	86	78
市场营销学	74	87	85	69	90	80	77	84	91	74	70
财务管理	68	75	70	84	73	60	76	81	88	68	75
统计学	55	91	68	73	84	81	70	69	94	62	71

4.3 用 pandas 调用 plot 函数画图

4.4 画子图函数 subplot



第 5 章 数据的概括性度量



5.1 均值、方差

5.2 众数、中位数、分位数

5.3 偏态、峰态

第 6 章 参数估计与假设检验



第 7 章 方差分析



第 8 章 线性回归



第 9 章 非线性回归



第 10 章 主成分分析



第 11 章 因子分析



第 12 章 时间序列分析和预测

