

Frequently Used Python Package

Requirements

Install Required Packages

```
1 | pip install -r requirements.txt
```

Fast Generated the requirements.txt

```
1 | pip freeze > requirements.txt
```

NumPy

NumPy (Numeric Python) 是 Python 中科学计算的基础包，它是一个 Python 库，提供多维数组对象，各种派生对象（如掩码数组和矩阵），以及用于数组快速操作的各种 API，有包括数学、逻辑、形状操作、排序、选择、输入输出、离散傅立叶变换、基本线性代数，基本统计运算和随机模拟等等

NumPy 包的核心是 `ndarray` 对象，它封装了 Python 原生的同数据类型的 `n` 维数组，为了保证其性能优良，其中有许多操作都是代码在本地进行编译后执行的，NumPy 相当于将 Python 变成一种免费的更强大的 MatLab 系统

NumPy 数组 和 原生 Python Array（数组）之间有几个重要的区别

- NumPy 数组在创建时具有固定的大小，与 Python 的原生数组对象（可以动态增长）不同，更改 `ndarray` 的大小将创建一个新数组并删除原来的数组
- NumPy 数组中的元素都需要具有相同的数据类型，因此在内存中的大小相同，例外情况，Python 的原生数组里包含了 NumPy 的对象的时候，这种情况下就允许不同大小元素的数组
- NumPy 数组有助于对大量数据进行高级数学和其他类型的操作，通常，这些操作的执行效率更高，比使用 Python 原生数组的代码更少
- 越来越多的基于 Python 的科学和数学软件包使用 NumPy 数组；虽然这些工具通常都支持 Python 的原生数组作为参数，但它们在处理之前会还是会将输入的数组转换为 NumPy 的数组，而且也通常输出为 NumPy 数组，为了高效地使用当今科学/数学基于 Python 的工具（大部分的科学计算工具），知道如何使用 NumPy 数组是必备的

扩展阅读 [【NumPy 中文文档】](#)

ndarray

From List to Array

列成 array-like 结构的数值数据可以通过使用 `array()` 函数转换为数组

```
1 In [1]: import numpy as np
2
3 In [2]: lst=[[1,2,3],[2,4,6]]
4
5 In [3]: type(lst)
6 Out[3]: list
7
8 In [4]: np_lst=np.array(lst)
9
10 In [5]: type(np_lst)
11 Out[5]: numpy.ndarray
```

Data Types of Arrays

```
1 In [6]: np_lst=np.array(lst, dtype=np.float)
2
3 In [7]: np_lst.dtype
4 Out[7]: dtype('float64')
```

Parameters of Array

```
1 In [8]: np_lst.shape      # 行列数
2 Out[8]: (2, 3)
3
4 In [9]: np_lst.ndim       # 维数
5 Out[9]: 2
6
7 In [10]: np_lst.itemsize  # 每个数据的数据存储大小
8 Out[10]: 8
9
10 In [11]: np_lst.size     # 元素个数
11 Out[11]: 6
```

Some kinds of Array

```
1 In [12]: np.zeros([2, 4])      # 生成 2 行 4 列的全 0 的数组
2 Out[12]:
3 array([[0., 0., 0., 0.],
4        [0., 0., 0., 0.]])
5
6 In [13]: np.ones([3, 5])      # 生成 3 行 5 列的全 1 的数组
7 Out[13]:
8 array([[1., 1., 1., 1., 1.],
9        [1., 1., 1., 1., 1.],
10       [1., 1., 1., 1., 1.]])
11
```

```

12 In [14]: np.random.rand(2, 4)    # 生成 2 行 4 列数组, 每个元素为 0-1 内均匀分布
      随机数
13 Out[14]:
14 array([[0.3492512 , 0.53278383, 0.67421472, 0.37741499],
15        [0.13505288, 0.56624554, 0.05743534, 0.47994088]])
16
17 In [15]: np.random.randint(1, 10, 3)    # 生成 3 个 1-10 内随机分布整数
18 Out[15]: array([6, 3, 5])
19
20 In [16]: np.random.randn(2, 4)          # 生成 2 行 4 列的标准正态随机数数组
21 Out[16]:
22 array([[ 0.7297433 , -1.31910919,  1.3258419 , -0.37062597],
23        [ 0.91714998,  2.0291667 ,  0.59648187, -1.54048607]])
24
25 In [17]: np.random.choice([10, 20, 30]) # 指定范围内的随机数
26 Out[17]: 20
27
28 In [18]: np.random.beta(1, 10, 20)      # 生成一个包含 20 个元素满足 Beta 分布
      的数组
29 Out[18]:
30 array([0.01745944, 0.19434248, 0.08223912, 0.04432289, 0.2939484 ,
31        0.13065389, 0.05528825, 0.20747935, 0.00320723, 0.11942977,
32        0.00388593, 0.00574769, 0.07600872, 0.08523846, 0.13702178,
33        0.01265392, 0.11381335, 0.01214367, 0.0733919 , 0.0779095 ])

```

Array Opeartion

Mathematical Operations of Array

```

1 In [19]: lst = np.arange(1, 11).reshape([2, 5])
2
3 In [20]: lst
4 Out[20]:
5 array([[ 1,  2,  3,  4,  5],
6        [ 6,  7,  8,  9, 10]])
7
8 In [21]: np.exp(lst)
9 Out[21]:
10 array([[2.71828183e+00, 7.38905610e+00, 2.00855369e+01, 5.45981500e+01,
11        1.48413159e+02],
12        [4.03428793e+02, 1.09663316e+03, 2.98095799e+03, 8.10308393e+03,
13        2.20264658e+04]])
14
15 In [22]: np.exp2(lst)
16 Out[22]:
17 array([[ 2.,  4.,  8., 16., 32.],
18        [ 64., 128., 256., 512., 1024.]])
19
20 In [23]: np.sqrt(lst)

```

```

19 Out[23]:
20 array([[1.          , 1.41421356, 1.73205081, 2.          , 2.23606798],
21        [2.44948974, 2.64575131, 2.82842712, 3.          , 3.16227766]])
22
23 In [24]: np.sin(lst)
24 Out[24]:
25 array([[ 0.84147098,  0.90929743,  0.14112001, -0.7568025 , -0.95892427],
26        [-0.2794155 ,  0.6569866 ,  0.98935825,  0.41211849, -0.54402111]])
27
28 In [25]: np.log(lst)
29 Out[25]:
30 array([[0.          , 0.69314718, 1.09861229, 1.38629436, 1.60943791],
31        [1.79175947, 1.94591015, 2.07944154, 2.19722458, 2.30258509]])

```

Descriptive Statistics for Array

```

1 In [26]: lst = np.array([[[1, 2, 3, 4], [4, 5, 6, 7]],
2     ...:                  [[7, 8, 9, 10], [10, 11, 12, 13]],
3     ...:                  [[14, 15, 16, 17], [18, 19, 20, 21]]])
4
5 In [27]: lst
6 Out[27]:
7 array([[[ 1,  2,  3,  4],
8         [ 4,  5,  6,  7]],
9
10        [[ 7,  8,  9, 10],
11         [10, 11, 12, 13]],
12
13        [[14, 15, 16, 17],
14         [18, 19, 20, 21]]])
15
16 In [28]: lst.sum()          # 所有元素求和
17 Out[28]: 252
18
19 In [29]: lst.sum(axis=0)    # 最外层求和
20 Out[29]:
21 array([[22, 25, 28, 31],
22        [32, 35, 38, 41]])
23
24 In [30]: lst.sum(axis=1)    # 第二层求和
25 Out[30]:
26 array([[ 5,  7,  9, 11],
27        [17, 19, 21, 23],
28        [32, 34, 36, 38]])
29
30 In [31]: lst.sum(axis=-1)   # 最里层求和
31 Out[31]:
32 array([[10, 22],
33        [34, 46],

```

```
34         [62, 78]])
35
36 In [32]: lst.max()
37 Out[32]: 21
38
39 In [33]: lst.min()
40 Out[33]: 1
```

Operations between Arrays

```
1 In [34]: lst1 = np.array([10, 20, 30, 40])
2
3 In [35]: lst2 = np.array([4, 3, 2, 1])
4
5 In [36]: lst1 + lst2
6 Out[36]: array([14, 23, 32, 41])
7
8 In [37]: lst1 - lst2
9 Out[37]: array([ 6, 17, 28, 39])
10
11 In [38]: lst1 * lst2
12 Out[38]: array([40, 60, 60, 40])
13
14 In [39]: lst1 / lst2
15 Out[39]: array([ 2.5,  6.66666667, 15.,  40.])
16
17 In [40]: lst1 ** lst2
18 Out[40]: array([10000,  8000,  900,  40])
19
20 In [41]: np.dot(lst1.reshape([2, 2]), lst2.reshape([2, 2]))
21 Out[41]:
22 array([[ 80,  50],
23        [200, 130]])
24
25 In [42]: np.concatenate((lst1, lst2), axis=0) # 向量拼接
26 Out[42]: array([10, 20, 30, 40,  4,  3,  2,  1])
27
28 In [43]: np.vstack((lst1, lst2)) # 按照行拼接
29 Out[43]:
30 array([[10, 20, 30, 40],
31        [ 4,  3,  2,  1]])
32
33 In [44]: np.hstack((lst1, lst2)) # 按照列拼接
34 Out[44]: array([10, 20, 30, 40,  4,  3,  2,  1])
35
36 In [45]: np.split(lst1, 2) # 向量拆分
37 Out[45]: [array([10, 20]), array([30, 40])]
38
39 In [46]: np.copy(lst1) # 向量拷贝
```

```
40 out[46]: array([10, 20, 30, 40])
```

Linear Algebra

```
1 In [47]: np.eye(3)      # 生成单位矩阵
2 Out[47]:
3 array([[1., 0., 0.],
4        [0., 1., 0.],
5        [0., 0., 1.]])
6
7 In [48]: from numpy.linalg import *
8
9 In [49]: lst = np.array([[1, 2],
10      ...:                [3, 4]])
11
12 In [50]: inv(lst)       # 生成给定矩阵的逆矩阵
13 Out[50]:
14 array([[ -2. ,  1. ],
15        [ 1.5, -0.5]])
16
17 In [51]: lst.transpose() # 生成给定矩阵的转置
18 Out[51]:
19 array([[1, 3],
20        [2, 4]])
21
22 In [52]: det(lst)       # 求矩阵的行列式
23 Out[52]: -2.0000000000000004
24
25 In [53]: eig(lst)       # 求矩阵的特征值和特征向量
26 Out[53]:
27 (array([-0.37228132,  5.37228132]),
28  array([[ -0.82456484, -0.41597356],
29        [ 0.56576746, -0.90937671]]))
```

注意

关于特征值和特征向量，例子中的两个特征值分别为

$$\lambda_1 = -0.37228132, \quad \lambda_2 = 5.37228132$$

对应的特征向量分别为

$$\xi_1 = \begin{bmatrix} -0.82456484 \\ 0.56576746 \end{bmatrix}, \quad \xi_2 = \begin{bmatrix} -0.41597356 \\ -0.90937671 \end{bmatrix}$$

可以验证如下等式

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} -0.82456484 \\ 0.56576746 \end{bmatrix} = -0.37228132 \times \begin{bmatrix} -0.82456484 \\ 0.56576746 \end{bmatrix}$$

Solve Linear Equations

```
1 In [54]: y = np.array([[5], [7]])
2
3 In [55]: solve(lst, y)
4 Out[55]:
5 array([[ -3.],
6         [ 4.]])
```

注意 相当于求解如下线性方程组

$$\begin{cases} x + 2y = 5 \\ 3x + 4y = 7 \end{cases}$$

NumPy Others

```
1 In [56]: np.corrcoef([1, 0, 1], [0, 2, 1])
2 Out[56]:
3 array([[ 1.          , -0.8660254],
4        [-0.8660254,  1.          ]])
5
6 In [57]: p = np.poly1d([2, 1, 3])      # 定义一元多项式 2 * x^2 + x + 3
7
8 In [58]: p(0.5)      # 多项式在 x = 0.5 时的值
9 Out[58]: 4.0
10
11 In [59]: p.r      # 多项式等于 0 时的根
12 Out[59]: array([-0.25+1.19895788j, -0.25-1.19895788j])
13
14 In [60]: q = np.poly1d([2, 1, 3], True)  # 把数组中的值作为根, 反推多项式
15
16 In [61]: print(q)
17      3      2
18 1 x - 6 x + 11 x - 6
```

注意 把数组中的值作为根, 反推多项式, 即

$$(x - 2)(x - 1)(x - 3) = x^3 - 6x^2 + 11x - 6$$

Matplotlib

如果要想象两个变量之间的关系, 想要显示值随时间变化, 就需要用到可视化工具

简单来说, Matplotlib 提供图形可视化 Python 包, 它提供了一种高度交互式界面, 便于用户能够做出各种有吸引力的统计图表

我们只需几行代码就可以生成图表、直方图、功率谱、条形图、误差图、散点图等

为了简单绘图，该 `pyplot` 模块提供了类似于 MATLAB 的界面，尤其是与 IPython 结合使用时。对于高级用户，您可以通过面向对象的界面或 MATLAB 用户熟悉的一组功能来完全控制线型，字体属性，轴属性等

扩展阅读 [【Matplotlib 中文文档】](#)

特别提醒，当使用 Jupyter 时，需要增加如下命令

```
1 | %matplotlib inline
```

具体作用是当你调用 `matplotlib.pyplot` 的绘图函数 `plot()` 等进行绘图的时候，或者生成一个 `figure` 画布的时候，可以直接在你的 Python Console 里面生成图像，但在 Python 的 IDE 如 Spyder 或者 Pycharm 中需要注释掉

Line

```
1 In [62]: import matplotlib.pyplot as plt
2
3 In [63]: def plt1():
4     ...:     x = np.linspace(-np.pi, np.pi, 256, endpoint=True)
5     ...:     c, s = np.cos(x), np.sin(x)
6     ...:     plt.figure(figsize=(10, 8))
7     ...:     # plt.plot(x, c)
8     ...:     plt.plot(x, c, color="blue", linewidth=1.5, linestyle="--",
9     ...:               label="cos", alpha=0.6) # 散点图
10    ...:     plt.plot(x, s, "r*", label="SIN", alpha=0.6)
11    ...:     plt.title("Cos & Sin", size=16) # 标题
12    ...:     ax = plt.gca() # 轴编辑器
13    ...:     ax.spines["right"].set_color("none")
14    ...:     ax.spines["top"].set_color("none")
15    ...:     ax.spines["left"].set_position(("data", 0))
16    ...:     ax.spines["bottom"].set_position(("data", 0))
17    ...:     ax.xaxis.set_ticks_position("bottom")
18    ...:     ax.yaxis.set_ticks_position("left")
19    ...:     plt.xticks([-np.pi, -np.pi / 2, 0, np.pi / 2, np.pi],
20    ...:               [r'$-\pi$', r'$-\pi/2$', r'$0$', r'$\pi/2$',
21    ...:               r'$\pi$']) # 正则表达
22    ...:     plt.yticks(np.linspace(-1, 1, 5, endpoint=True))
23    ...:     for label in ax.get_xticklabels() + ax.get_yticklabels():
24    ...:         label.set_fontsize(16)
25    ...:         label.set_bbox(dict(facecolor="white", edgecolor="none",
26    ...:                               alpha=0.2))
27    ...:     plt.legend(loc="upper left") # 图例位置
28    ...:     plt.grid() # 网格线
29    ...:     # fill
30    ...:     plt.fill_between(x, np.abs(x) < 0.5, c, c > 0.5,
31    ...:                     color="green", alpha=0.25)
32    ...:     t = 1
```



```

31     ...:     plt.plot([t, t], [0, np.cos(t)], "y", linewidth=3,
32         linestyle="--") # 橙色虚线
33     ...:     plt.annotate("cos(1)", xy=(t, np.cos(1)), xycoords="data",
34         xytext=(+10, +30), textcoords="offset points",
35         arrowprops=dict(arrowstyle="->",
36         connectionstyle="arc3,rad=.2"))
37     ...:     plt.show() # 显示
38 In [64]: plt1()

```

在填充画图代码中

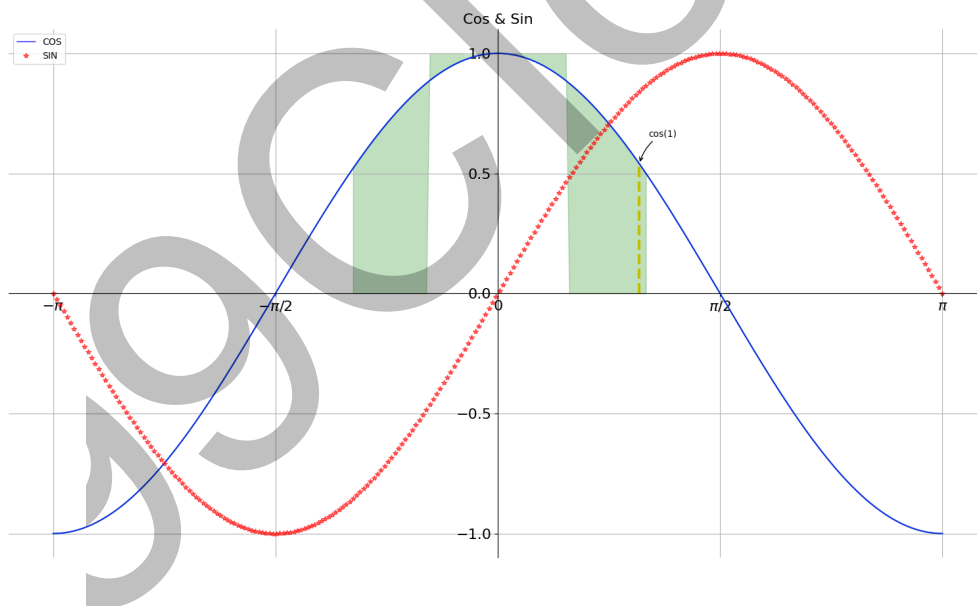
```

1 plt.fill_between(x, np.abs(x) < 0.5, c, c > 0.5, color="green", alpha=0.25)

```

第一个参数 `x` 表示 x 轴，第二个参数 `np.abs(x)` 表示 x 的绝对值，`np.abs(x) < 0.5` 是一个判定变量，`c` 表示 y 轴，`c > 0.5` 是一个判定条件

- 当 `np.abs(x) < 0.5` 为 `True`（即值为 1），从 y 轴的 1（满足 $c > 0.5$ ）开始往两边填充（当然 x 轴上是 -0.5 到 0.5 之间的区域），此时填充的也就是图上方的两小块
- 当 `np.abs(x) >= 0.5` 为 `False`（即值为 0），从 y 轴的 0 开始向上填充，当然只填充 $c > 0.5$ 的区域，也就是图中那两块大的对称区域



Style

```

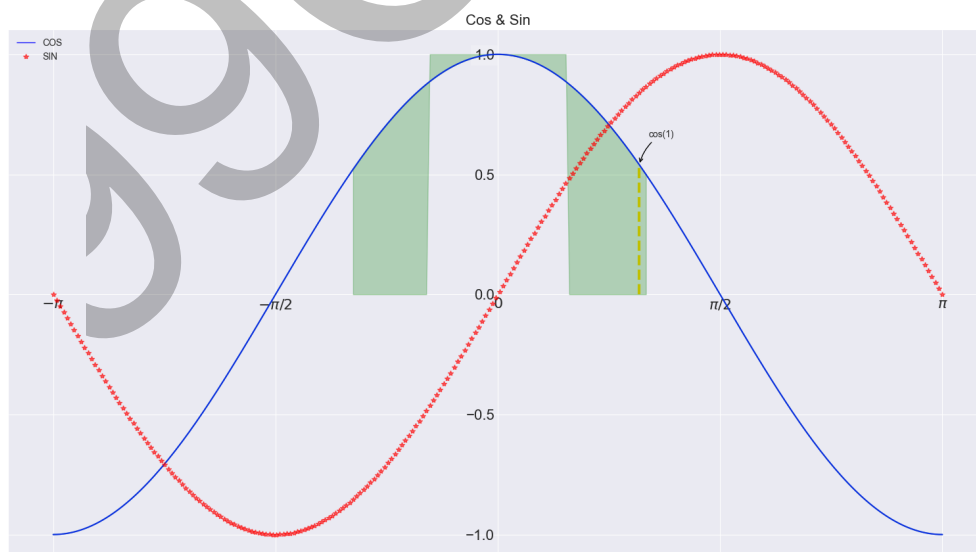
1 In [65]: plt.style.available # 查看可用画风
2 Out[65]:
3 ['Solarize_Light2',
4  '_classic_test_patch',
5  'bmh',
6  'classic',

```

```

7 'dark_background',
8 'fast',
9 'fivethirtyeight',
10 'ggplot',
11 'grayscale',
12 'seaborn',
13 'seaborn-bright',
14 'seaborn-colorblind',
15 'seaborn-dark',
16 'seaborn-dark-palette',
17 'seaborn-darkgrid',
18 'seaborn-deep',
19 'seaborn-muted',
20 'seaborn-notebook',
21 'seaborn-paper',
22 'seaborn-pastel',
23 'seaborn-poster',
24 'seaborn-talk',
25 'seaborn-ticks',
26 'seaborn-white',
27 'seaborn-whitegrid',
28 'tableau-colorblind10']
29
30 In [66]: plt.style.use('seaborn-dark') # 应用风格
31
32 In [67]: plt1()
33
34 In [68]: plt.style.use('default') # 重回默认风格

```



Many types of Figures

```

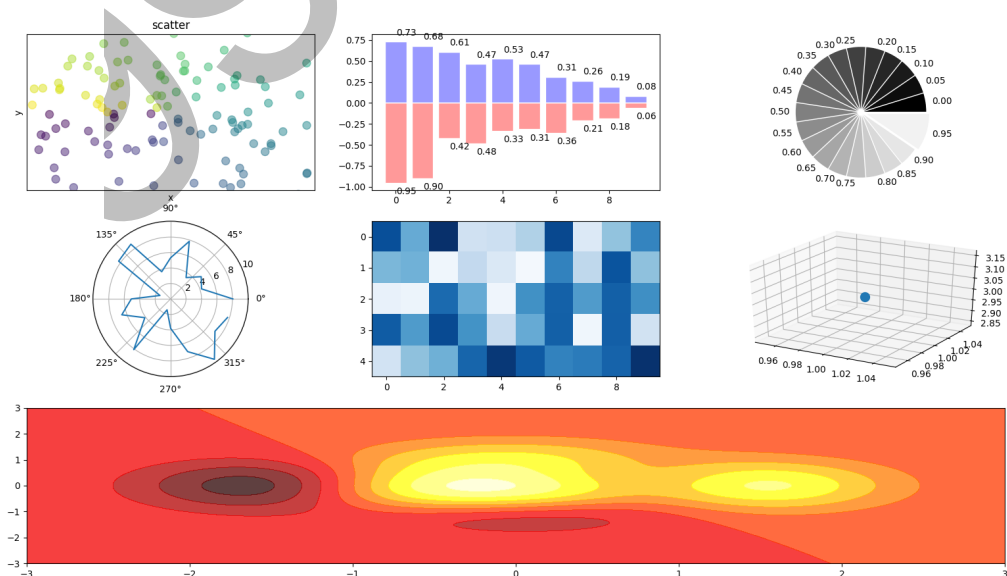
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4
5 def plt2():
6     fig = plt.figure()
7     # scatter
8     ax = fig.add_subplot(3, 3, 1)
9     n = 128
10    X = np.random.normal(0, 1, n)
11    Y = np.random.normal(0, 1, n)
12    T = np.arctan2(Y, X)
13    # plt.axes([0.025, 0.025, 0.95, 0.95])
14    plt.scatter(X, Y, s=75, c=T, alpha=.5)
15    plt.xlim(-1.5, 1.5), plt.xticks([])
16    plt.ylim(-1.5, 1.5), plt.yticks([])
17    plt.axis()
18    plt.title("scatter")
19    plt.xlabel("x")
20    plt.ylabel("y")
21
22    # bar
23    fig.add_subplot(332)
24    n = 10
25    X = np.arange(n)
26    Y1 = (1 - X / float(n)) * np.random.uniform(0.5, 1, n)
27    Y2 = (1 - X / float(n)) * np.random.uniform(0.5, 1, n)
28    plt.bar(X, +Y1, facecolor='#9999ff', edgecolor='white')
29    plt.bar(X, -Y2, facecolor='#ff9999', edgecolor='white')
30    for x, y in zip(X, Y1):
31        plt.text(x + 0.4, y + 0.05, '%.2f' % y, ha='center', va='bottom')
32    for x, y in zip(X, Y2):
33        plt.text(x + 0.4, -y - 0.05, '%.2f' % y, ha='center', va='top')
34
35    # Pie
36    fig.add_subplot(333)
37    n = 20
38    Z = np.ones(n)
39    Z[-1] *= 2
40    # explode扇形离中心距离
41    plt.pie(Z, explode=Z * .05, colors=['%f' % (i / float(n)) for i in
range(n)],
42          labels=['%.2f' % (i / float(n)) for i in range(n)])
43    plt.gca().set_aspect('equal') # 圆形
44    plt.xticks([], plt.yticks([])
45
46    # polar
47    fig.add_subplot(334, polar=True)
48    n = 20

```

```

49 theta = np.arange(0, 2 * np.pi, 2 * np.pi / n)
50 radii = 10 * np.random.rand(n)
51 plt.polar(theta, radii)
52 # plt.plot(theta, radii)
53
54 # heatmap
55 fig.add_subplot(335)
56 from matplotlib import cm
57 data = np.random.rand(5, 10)
58 cmap = cm.Blues
59 map = plt.imshow(data, interpolation='nearest', cmap=cmap,
60 aspect='auto', vmin=0, vmax=1)
61
62 # 3D
63 from mpl_toolkits.mplot3d import Axes3D
64 ax = fig.add_subplot(336, projection="3d")
65 ax.scatter(1, 1, 3, s=100)
66
67 # hot map
68 fig.add_subplot(313)
69
70 def f(x, y):
71     return (1 - x / 2 + x ** 5 + y ** 3) * np.exp(- x ** 2 - y ** 2)
72
73 n = 256
74 x = np.linspace(-3, 3, n * 2)
75 y = np.linspace(-3, 3, n)
76 X, Y = np.meshgrid(x, y)
77 plt.contourf(X, Y, f(X, Y), 8, alpha=.75, cmap=plt.cm.hot)
78
79 plt.show() # 显示

```



Scipy

SciPy 函数库在 NumPy 库的基础上增加了众多的数学、科学以及工程计算中常用的库函数，例如线性代数、常微分方程数值求解、信号处理、图像处理、稀疏矩阵等等

扩展阅读【[SciPy_官网](#)】

Solving Nonlinear Equations

optimize 库中的 `fsolve` 函数可以用来对非线性方程组进行求解，它的基本调用形式如下

```
1 fsolve(func, x0)
```

`func(x)` 是计算方程组误差的函数，它的参数 x 是一个矢量，表示方程组的各个未知数的一组可能解，`func` 返回将 x 代入方程组之后得到的误差； x_0 为未知数矢量的初始值，如果要对如下方程组进行求解的话

$$\begin{cases} f_1(u_1, u_2, u_3) = 0 \\ f_2(u_1, u_2, u_3) = 0 \\ f_3(u_1, u_2, u_3) = 0 \end{cases}$$

那么 `func` 可以如下定义

```
1 def func(x):
2     u1,u2,u3 = x
3     return [f1(u1,u2,u3), f2(u1,u2,u3), f3(u1,u2,u3)]
```

下面是一个实际的例子，求解如下方程组的解

$$\begin{cases} 5 \cdot x_1 + 3 = 0 \\ 4 \cdot x_0^2 - 2 \sin(x_1 \cdot x_2) = 0 \\ x_1 \cdot x_2 - 1.5 = 0 \end{cases}$$

程序如下

```
1 from scipy.optimize import fsolve
2 from math import sin
3
4
5 def f(x):
6     x0 = float(x[0])
7     x1 = float(x[1])
8     x2 = float(x[2])
9     return [
10         5 * x1 + 3,
11         4 * x0 * x0 - 2 * sin(x1 * x2),
12         x1 * x2 - 1.5
13     ]
```

```

14
15
16 result = fsolve(f, [1, 1, 1])
17
18 print('[x0,x1,x2] =', result)
19 print('[f1,f2,f3] =', f(result))

```

输出为

```

1 [x0,x1,x2] = [-0.70622057 -0.6          -2.5          ]
2 [f1,f2,f3] = [0.0, -9.126033262418787e-14, 5.329070518200751e-15]

```

由于 `fsolve` 函数在调用函数 `f` 时，传递的参数为数组，因此如果直接使用数组中的元素计算的话，计算速度将会有所降低，因此这里先用 `float` 函数将数组中的元素转换为 Python 中的标准浮点数，然后调用标准 `math` 库中的函数进行运算

Function Maximum

以寻找函数

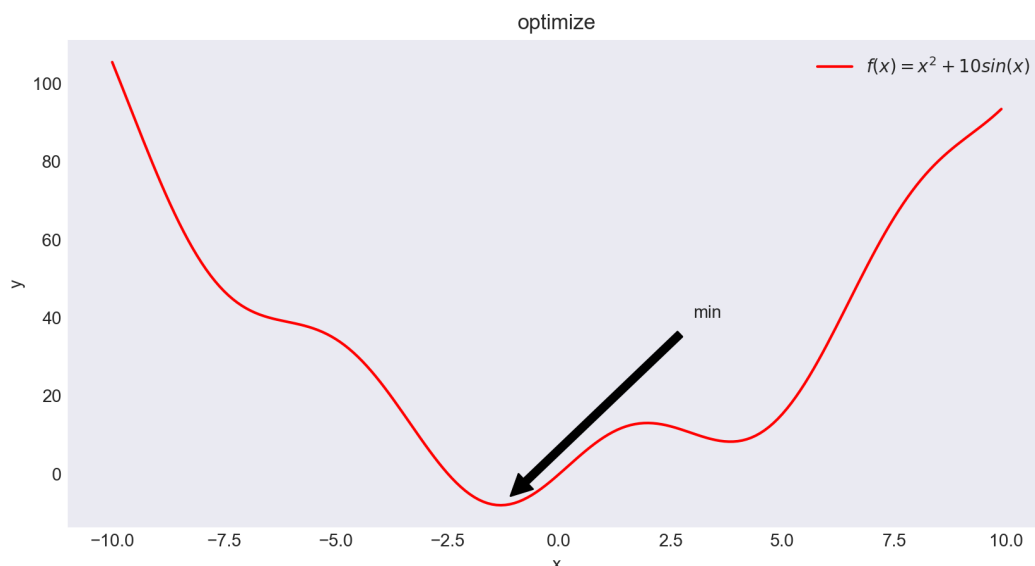
$$f(x) = x^2 + 10\sin(x)$$

的最小值为例，首先绘制目标函数的图形

```

1 from scipy import optimize
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5
6 # 定义目标函数
7 def f(x):
8     return x ** 2 + 10 * np.sin(x)
9
10
11 # 绘制目标函数的图形
12 plt.figure(figsize=(10, 5))
13 x = np.arange(-10, 10, 0.1)
14 plt.xlabel('x')
15 plt.ylabel('y')
16 plt.title('optimize')
17 plt.plot(x, f(x), 'r-', label='$f(x)=x^2+10sin(x)$')
18 # 图像中的最低点函数值
19 a = f(-1.3)
20 plt.annotate('min', xy=(-1.3, a), xytext=(3, 40),
21             arrowprops=dict(facecolor='black', shrink=0.05))
22 plt.legend()
23 plt.show()

```



显然这是一个非凸优化问题，对于这类函数得最小值问题一般是从给定的初始值开始进行一个梯度下降，在 `optimize` 中一般使用 `bfgs` 算法

```
1 optimize.fmin_bfgs(f, 0)
```

结果显示在经过五次迭代之后找到了一个局部最低点 `-7.945823`，显然这并不是函数的全局最小值，只是该函数的一个局部最小值，这也是拟牛顿算法（BFGS）的局限性，如果一个函数有多个局部最小值，拟牛顿算法可能找到这些局部最小值而不是全局最小值，这取决于初始点的选取

```
1 Optimization terminated successfully.
2     Current function value: -7.945823
3     Iterations: 5
4     Function evaluations: 18
5     Gradient evaluations: 6
```

在我们不知道全局最低点，并且使用一些临近点作为初始点，那将需要花费大量的时间来获得全局最优，此时可以采用暴力搜寻算法，它会评估范围网格内的每一个点，对于本例，如下

```
1 grid = (-10, 10, 0.1)
2 xmin_global = optimize.brute(f, (grid,))
3 print(xmin_global)
```

搜寻结果如下

```
1 [-1.30641113]
```

但是当函数的定义域大到一定程度时，`scipy.optimize.brute()` 变得非常慢，`scipy.optimize.basinhopping()` 提供了一个解决思路

```

1 x0 = -10
2 xmin_global_2 = optimize.basinhopping(f, x0, stepsize=5).x
3 print(xmin_global_2)

```

搜寻结果如下

```

1 [-1.30644]

```

Least Square Fitting

假设有一组实验数据 $(x[i], y[i])$ ，我们知道它们之间的函数关系 $y = f(x)$ ，通过这些已知信息，需要确定函数中的一些参数项

例如，如果 f 是一个线型函数 $f(x) = k \times x + b$ ，那么参数 k 和 b 就是我们需要确定的值，如果将这些参数用 \mathbf{p} 表示的话，那么我们就是要找到一组 \mathbf{p} 值使得如下公式中的 S 函数最小

$$S(\mathbf{p}) = \sum_{i=1}^m [y_i - f(x_i, \mathbf{p})]^2$$

这种算法被称之为最小二乘拟合 (Least-square fitting)

scipy 中的子函数库 optimize 已经提供了实现最小二乘拟合算法的函数 `leastsq`

下面是用 `leastsq` 进行数据拟合的一个例子

```

1 import numpy as np
2 from scipy.optimize import leastsq
3 import pylab as pl
4 pl.mpl.rcParams['font.sans-serif'] = ['SimHei']
5 pl.mpl.rcParams['axes.unicode_minus'] = False
6
7
8 def func(x, p):
9     """
10     数据拟合所用的函数: A*sin(2*pi*k*x + theta)
11     """
12     A, k, theta = p
13     return A*np.sin(2*np.pi*k*x+theta)
14
15
16 def residuals(p, y, x):
17     """
18     实验数据x, y和拟合函数之间的差, p为拟合需要找到的系数
19     """
20     return y - func(x, p)
21
22
23 x = np.linspace(0, - 2 * np.pi, 100)
24 A, k, theta = 10, 0.34, np.pi/6 # 真实数据的函数参数

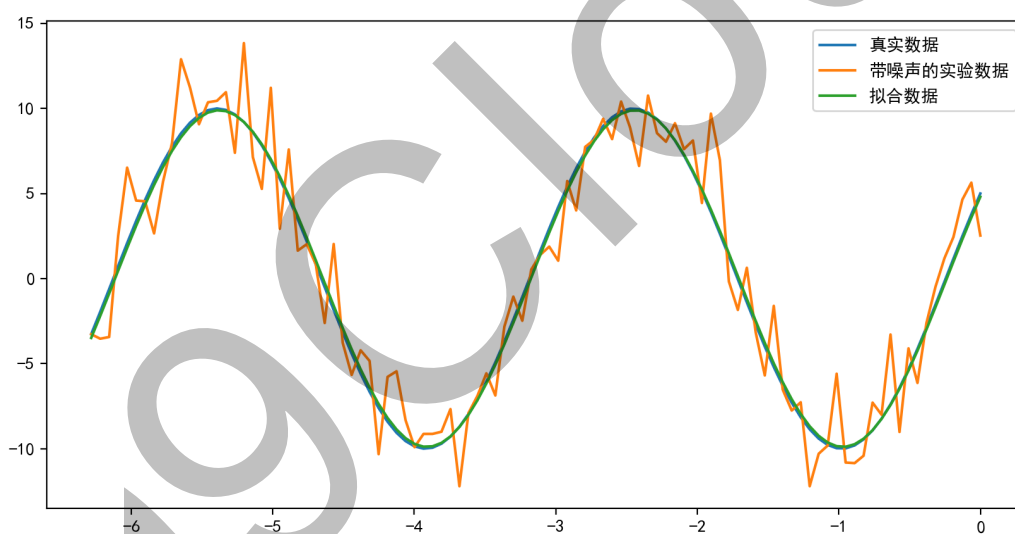
```



```

25 y0 = func(x, [A, k, theta])          # 真实数据
26 y1 = y0 + 2 * np.random.randn(len(x)) # 加入噪声之后的实验数据
27
28 p0 = [7, 0.2, 0]                    # 第一次猜测的函数拟合参数
29
30 # 调用 leastsq 进行数据拟合
31 # residuals 为计算误差的函数
32 # p0 为拟合参数的初始值
33 # args 为需要拟合的实验数据
34 plsq = leastsq(residuals, p0, args=(y1, x))
35
36 print(u"真实参数:", [A, k, theta])
37 print(u"拟合参数:", plsq[0]) # 实验数据拟合后的参数
38
39 pl.plot(x, y0, label=u"真实数据")
40 pl.plot(x, y1, label=u"带噪声的实验数据")
41 pl.plot(x, func(x, plsq[0]), label=u"拟合数据")
42 pl.legend()
43 pl.show()

```



输出结果

```

1  真实参数: [10, 0.34, 0.5235987755982988]
2  拟合参数: [10.22216161  0.34359989  0.50580946]

```

这个例子中我们要拟合的函数是一个正弦波函数，它有三个参数 **A**, **k**, **theta**，分别对应振幅、频率、相角，假设我们的实验数据是一组包含噪声的数据 x, y_1 ，其中 y_1 是在真实数据 y_0 的基础上加入噪声得到的

Pandas


```

42     ...:                                     "C": pd.Series(1, index=list(range(4)),
dtype="float32"),
43     ...:                                     "D": np.array([3] * 4, dtype="float32"),
44     ...:                                     "E": pd.Categorical(["police", "student",
"teacher", "doctor"])}))
45
46 In [10]: df
47 Out[10]:
48      A      B      C      D      E
49 0  1 2020-05-01  1.0  3.0  police
50 1  1 2020-05-01  1.0  3.0  student
51 2  1 2020-05-01  1.0  3.0  teacher
52 3  1 2020-05-01  1.0  3.0  doctor

```

Basic & Select & Set

我们先重新设置一个 `df`

```

1 In [11]: dates = pd.date_range("20200501", periods=8)
2
3 In [12]: df = pd.DataFrame(np.random.randn(8, 5), index=dates,
columns=list("ABCDE"))
4
5 In [13]: df
6 Out[13]:
7      A      B      C      D      E
8 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
9 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
10 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396
11 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240
12 2020-05-05  0.343573  2.268671  0.123661  0.026515  0.819831
13 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408
14 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575
15 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486

```

Basic Operation of DataFrame

```

1 In [14]: df.head(3)
2 Out[14]:
3      A      B      C      D      E
4 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
5 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
6 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396
7
8 In [15]: df.tail(3)
9 Out[15]:
10      A      B      C      D      E
11 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408

```

```

12 2020-05-07 -0.504259 0.648047 -1.170642 1.198958 0.756575
13 2020-05-08 -0.521412 0.176260 -0.090847 0.020373 0.047486
14
15 In [16]: df.index
16 Out[16]:
17 DatetimeIndex(['2020-05-01', '2020-05-02', '2020-05-03', '2020-05-04',
18               '2020-05-05', '2020-05-06', '2020-05-07', '2020-05-08'],
19               dtype='datetime64[ns]', freq='D')
20
21 In [17]: df.values
22 Out[17]:
23 array([[ -0.63960556, -0.49076281,  0.83405679,  0.19167819,  0.54413525],
24        [  0.74831337,  0.54242504, -0.50424864, -1.00617742,  1.08873712],
25        [  0.56241309, -1.36623581, -0.9994565 ,  1.43409736,  0.08839574],
26        [-0.01196919,  1.60814395,  1.0709043 ,  1.89549629,  1.37623995],
27        [  0.34357293,  2.26867082,  0.12366068,  0.02651487,  0.8198315 ],
28        [  1.14905497,  0.51412852,  0.65329835,  1.01494828,  0.36440802],
29        [-0.50425931,  0.64804705, -1.17064214,  1.19895757,  0.7565752 ],
30        [-0.52141233,  0.1762601 , -0.09084731,  0.0203732 ,  0.04748562]])
31
32 In [18]: df.T
33 Out[18]:
34      2020-05-01  2020-05-02  2020-05-03  ...  2020-05-06  2020-05-07  2020-
35      05-08
36  A    -0.639606    0.748313    0.562413  ...    1.149055   -0.504259
37      -0.521412
38  B    -0.490763    0.542425   -1.366236  ...    0.514129    0.648047
39      0.176260
40  C     0.834057   -0.504249   -0.999457  ...    0.653298   -1.170642
41      -0.090847
42  D     0.191678   -1.006177    1.434097  ...    1.014948    1.198958
43      0.020373
44  E     0.544135    1.088737    0.088396  ...    0.364408    0.756575
45      0.047486
46
47 [5 rows x 8 columns]
48
49 In [19]: df.sort_values(by="C", ascending=False)
50 Out[19]:
51           A         B         C         D         E
52 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240
53 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
54 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408
55 2020-05-05  0.343573  2.268671  0.123661  0.026515  0.819831
56 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486
57 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
58 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396
59 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575

```

```

55 In [20]: df.sort_index(axis=1, ascending=False)
56 Out[20]:
57          E          D          C          B          A
58 2020-05-01  0.544135  0.191678  0.834057 -0.490763 -0.639606
59 2020-05-02  1.088737 -1.006177 -0.504249  0.542425  0.748313
60 2020-05-03  0.088396  1.434097 -0.999457 -1.366236  0.562413
61 2020-05-04  1.376240  1.895496  1.070904  1.608144 -0.011969
62 2020-05-05  0.819831  0.026515  0.123661  2.268671  0.343573
63 2020-05-06  0.364408  1.014948  0.653298  0.514129  1.149055
64 2020-05-07  0.756575  1.198958 -1.170642  0.648047 -0.504259
65 2020-05-08  0.047486  0.020373 -0.090847  0.176260 -0.521412
66
67 In [21]: df.describe()      # 描述统计
68 Out[21]:
69          A          B          C          D          E
70 count  8.000000  8.000000  8.000000  8.000000  8.000000
71 mean   0.140763  0.487585 -0.010409  0.596986  0.635726
72 std    0.664565  1.130620  0.837986  0.949908  0.467467
73 min   -0.639606 -1.366236 -1.170642 -1.006177  0.047486
74 25%   -0.508548  0.009504 -0.628051  0.024979  0.295405
75 50%    0.165802  0.528277  0.016407  0.603313  0.650355
76 75%    0.608888  0.888071  0.698488  1.257743  0.887058
77 max    1.149055  2.268671  1.070904  1.895496  1.376240

```

Selection Operations of DataFrame

```

1 In [22]: type(df["A"])
2 Out[22]: pandas.core.series.Series
3
4 In [23]: df[:3]
5 Out[23]:
6          A          B          C          D          E
7 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
8 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
9 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396
10
11 In [24]: df.head(3)
12 Out[24]:
13          A          B          C          D          E
14 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
15 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
16 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396
17
18 In [25]: df["20200501": "20200504"]
19 Out[25]:
20          A          B          C          D          E
21 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135
22 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737
23 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396

```

```

24 2020-05-04 -0.011969 1.608144 1.070904 1.895496 1.376240
25
26 In [26]: df.loc["20200501": "20200504", ["B", "D"]]      # 行名和列名
27 Out[26]:
28              B          D
29 2020-05-01 -0.490763  0.191678
30 2020-05-02  0.542425 -1.006177
31 2020-05-03 -1.366236  1.434097
32 2020-05-04  1.608144  1.895496
33
34 In [27]: df.at[dates[0], "C"]
35 Out[27]: 0.8340567905719413
36
37 In [28]: df.iloc[1:3, 2:4]                                # 行号和列号
38 Out[28]:
39              C          D
40 2020-05-02 -0.504249 -1.006177
41 2020-05-03 -0.999457  1.434097
42
43 In [29]: df.iloc[1, 4]
44 Out[29]: 1.0887371182725154
45
46 In [30]: df.iat[1, 4]
47 Out[30]: 1.0887371182725154

```

Data Intercepted by Judgment in DataFrame

```

1 In [31]: df[df.B > 0][df.A < 0]
2 ...../anaconda3/envs/AIC/bin/ipython:1: UserWarning: Boolean Series key
   will be reindexed to match DataFrame index.
3   #!...../anaconda3/envs/AIC/bin/python
4 Out[31]:
5              A          B          C          D          E
6 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240
7 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575
8 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486
9
10 In [32]: df[(df.B > 0) & (df.A < 0)]
11 Out[32]:
12              A          B          C          D          E
13 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240
14 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575
15 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486

```

布尔型系列键将索引匹配获得对应的索引

我每一个判断之后，都会返回一个 `True` 和 `False` 的索引列表（矩阵），通过对数据索引位置的布尔判断来筛选条件，如果一条语句出现两个判断条件，会存在语义不明的情况

```

1 In [33]: df[df > 0]
2 Out[33]:
3           A          B          C          D          E
4 2020-05-01      NaN      NaN  0.834057  0.191678  0.544135
5 2020-05-02  0.748313  0.542425      NaN      NaN  1.088737
6 2020-05-03  0.562413      NaN      NaN  1.434097  0.088396
7 2020-05-04      NaN  1.608144  1.070904  1.895496  1.376240
8 2020-05-05  0.343573  2.268671  0.123661  0.026515  0.819831
9 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408
10 2020-05-07      NaN  0.648047      NaN  1.198958  0.756575
11 2020-05-08      NaN  0.176260      NaN  0.020373  0.047486
12
13 In [34]: df[df["E"].isin([1, 2])]
14 Out[34]:
15 Empty DataFrame
16 Columns: [A, B, C, D, E]
17 Index: []

```

Setting Value Operation in DataFrame

```

1 In [35]: s1 = pd.Series(list(range(10,18)),index=pd.date_range("20200501",
2 periods=8))
3
4 In [36]: df["F"] = s1
5
6 In [37]: df
7 Out[37]:
8           A          B          C          D          E          F
9 2020-05-01 -0.639606 -0.490763  0.834057  0.191678  0.544135  10
10 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737  11
11 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396  12
12 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240  13
13 2020-05-05  0.343573  2.268671  0.123661  0.026515  0.819831  14
14 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408  15
15 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575  16
16 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486  17
17
18 In [38]: df.at[dates[0], "A"] = 0
19
20 In [39]: df
21 Out[39]:
22           A          B          C          D          E          F
23 2020-05-01  0.000000 -0.490763  0.834057  0.191678  0.544135  10
24 2020-05-02  0.748313  0.542425 -0.504249 -1.006177  1.088737  11
25 2020-05-03  0.562413 -1.366236 -0.999457  1.434097  0.088396  12
26 2020-05-04 -0.011969  1.608144  1.070904  1.895496  1.376240  13
27 2020-05-05  0.343573  2.268671  0.123661  0.026515  0.819831  14
28 2020-05-06  1.149055  0.514129  0.653298  1.014948  0.364408  15
29 2020-05-07 -0.504259  0.648047 -1.170642  1.198958  0.756575  16

```

```

29 2020-05-08 -0.521412  0.176260 -0.090847  0.020373  0.047486  17
30
31 In [40]: df.iat[1, 1] = 1
32
33 In [41]: df.loc[:, "D"] = np.array([4] * len(df))
34
35 In [42]: df
36 Out[42]:
37           A          B          C  D          E  F
38 2020-05-01  0.000000 -0.490763  0.834057  4  0.544135 10
39 2020-05-02  0.748313  1.000000 -0.504249  4  1.088737 11
40 2020-05-03  0.562413 -1.366236 -0.999457  4  0.088396 12
41 2020-05-04 -0.011969  1.608144  1.070904  4  1.376240 13
42 2020-05-05  0.343573  2.268671  0.123661  4  0.819831 14
43 2020-05-06  1.149055  0.514129  0.653298  4  0.364408 15
44 2020-05-07 -0.504259  0.648047 -1.170642  4  0.756575 16
45 2020-05-08 -0.521412  0.176260 -0.090847  4  0.047486 17
46
47 In [43]: df2 = df.copy()
48
49 In [44]: df2[df2 > 0] = - df2
50
51 In [45]: df2
52           A          B          C  D          E  F
53 2020-05-01  0.000000 -0.490763 -0.834057 -4 -0.544135 -10
54 2020-05-02 -0.748313 -1.000000 -0.504249 -4 -1.088737 -11
55 2020-05-03 -0.562413 -1.366236 -0.999457 -4 -0.088396 -12
56 2020-05-04 -0.011969 -1.608144 -1.070904 -4 -1.376240 -13
57 2020-05-05 -0.343573 -2.268671 -0.123661 -4 -0.819831 -14
58 2020-05-06 -1.149055 -0.514129 -0.653298 -4 -0.364408 -15
59 2020-05-07 -0.504259 -0.648047 -1.170642 -4 -0.756575 -16
60 2020-05-08 -0.521412 -0.176260 -0.090847 -4 -0.047486 -17

```

Missing Data Processing

我们先重新设置一个 df

```

1 In [46]: dates = pd.date_range("20200501", periods=8)
2
3 In [47]: df = pd.DataFrame(np.random.randn(8, 5), index=dates,
4                             columns=list("ABCDE"))
5
6 In [48]: df
7 Out[48]:
8           A          B          C          D          E
9 2020-05-01  1.161780  0.046974  0.317034  0.985277 -0.878156
10 2020-05-02 -0.511518 -0.462444 -0.090256  1.013958 -0.052817
11 2020-05-03  0.492906 -0.098113 -1.621421 -0.469094 -0.954550
12 2020-05-04  0.229874 -0.344795 -0.158310 -0.419449  0.096488

```



```

12 2020-05-05  1.250265 -1.422900  1.084396  0.902803 -1.138471
13 2020-05-06 -1.349342 -0.357210 -0.623589  0.331251  0.305456
14 2020-05-07  0.506861 -1.480997 -0.835471  0.158394  1.484623
15 2020-05-08 -0.387560 -0.233622  1.192566 -0.510911 -0.855755

```

Missing Values

```

1 In [49]: df1 = df.reindex(index=dates[:4],
2     ...:                  columns=list("ABCD") + ["G"])
3
4 In [50]: df1.loc[dates[0]: dates[1], "G"] = 1
5
6 In [51]: df1
7 Out[51]:
8              A          B          C          D          G
9 2020-05-01  1.161780  0.046974  0.317034  0.985277  1.0
10 2020-05-02 -0.511518 -0.462444 -0.090256  1.013958  1.0
11 2020-05-03  0.492906 -0.098113 -1.621421 -0.469094  NaN
12 2020-05-04  0.229874 -0.344795 -0.158310 -0.419449  NaN
13
14 In [52]: df1.dropna()
15 Out[52]:
16              A          B          C          D          G
17 2020-05-01  1.161780  0.046974  0.317034  0.985277  1.0
18 2020-05-02 -0.511518 -0.462444 -0.090256  1.013958  1.0
19
20 In [53]: df1.fillna(value=2)
21 Out[53]:
22              A          B          C          D          G
23 2020-05-01  1.161780  0.046974  0.317034  0.985277  1.0
24 2020-05-02 -0.511518 -0.462444 -0.090256  1.013958  1.0
25 2020-05-03  0.492906 -0.098113 -1.621421 -0.469094  2.0
26 2020-05-04  0.229874 -0.344795 -0.158310 -0.419449  2.0

```

Merge & Reshape

Pandas Statistic

```

1 In [54]: df.mean()
2 Out[54]:
3 A    0.174158
4 B   -0.544138
5 C   -0.091881
6 D    0.249029
7 E   -0.249148
8 dtype: float64
9
10 In [55]: df.var()

```

```
11 Out[55]:
12 A      0.779362
13 B      0.339436
14 C      0.911456
15 D      0.444235
16 E      0.789096
17 dtype: float64
18
19 In [56]: s = pd.Series([1, 2, 2, np.nan, 5, 7, 9, 10], index=dates)
20
21 In [57]: s
22 Out[57]:
23 2020-05-01      1.0
24 2020-05-02      2.0
25 2020-05-03      2.0
26 2020-05-04      NaN
27 2020-05-05      5.0
28 2020-05-06      7.0
29 2020-05-07      9.0
30 2020-05-08     10.0
31 Freq: D, dtype: float64
32
33 In [58]: s.shift(2)          # shift 函数是对数据进行移动的操作
34 Out[58]:
35 2020-05-01      NaN
36 2020-05-02      NaN
37 2020-05-03      1.0
38 2020-05-04      2.0
39 2020-05-05      2.0
40 2020-05-06      NaN
41 2020-05-07      5.0
42 2020-05-08      7.0
43 Freq: D, dtype: float64
44
45 In [59]: s.diff()          # 差分列
46 Out[59]:
47 2020-05-01      NaN
48 2020-05-02      1.0
49 2020-05-03      0.0
50 2020-05-04      NaN
51 2020-05-05      NaN
52 2020-05-06      2.0
53 2020-05-07      2.0
54 2020-05-08      1.0
55 Freq: D, dtype: float64
56
57 In [60]: s.value_counts()   # 频数统计
58 Out[60]:
59 2.0      2
```

```

60 10.0    1
61  9.0    1
62  7.0    1
63  5.0    1
64  1.0    1
65 dtype: int64

```

Pandas Concat

```

1  In [61]: pieces = [df[:3], df[-3:]]
2
3  In [62]: pd.concat(pieces)
4  Out[62]:
5
6           A          B          C          D          E
7 2020-05-01  1.161780  0.046974  0.317034  0.985277 -0.878156
8 2020-05-02 -0.511518 -0.462444 -0.090256  1.013958 -0.052817
9 2020-05-03  0.492906 -0.098113 -1.621421 -0.469094 -0.954550
10 2020-05-06 -1.349342 -0.357210 -0.623589  0.331251  0.305456
11 2020-05-07  0.506861 -1.480997 -0.835471  0.158394  1.484623
12 2020-05-08 -0.387560 -0.233622  1.192566 -0.510911 -0.855755
13
14 In [63]: left = pd.DataFrame({"key": ["x", "y"], "value": [1, 2]})
15
16 In [64]: left
17 Out[64]:
18    key  value
19 0    x      1
20 1    y      2
21
22 In [65]: right = pd.DataFrame({"key": ["x", "z"], "value": [3, 4]})
23
24 In [66]: right
25 Out[66]:
26    key  value
27 0    x      3
28 1    z      4
29
30 In [67]: pd.merge(left, right, on="key", how="outer")
31 Out[67]:
32    key  value_x  value_y
33 0    x      1.0      3.0
34 1    y      2.0      NaN
35 2    z      NaN      4.0
36
37 In [68]: df3 = pd.DataFrame({"A": ["a", "b", "c", "b"], "B":
38    list(range(4))})
39
40 In [69]: df3.groupby("A").sum()      # a: 0; b: 1+3; c: 2
41 Out[69]:

```

```

40      B
41      A
42      a  0
43      b  4
44      c  2

```

Pandas Reshape

```

1  In [70]: import datetime
2
3  In [71]: df4 = pd.DataFrame({'A': ['one', 'one', 'two', 'three'] * 6,
4      ....:                   'B': ['A', 'B', 'C'] * 8,
5      ....:                   'C': ['foo', 'foo', 'foo', 'bar', 'bar',
6      ....:                   'bar'] * 4,
7      ....:                   'D': np.random.randn(24),
8      ....:                   'E': np.random.randn(24),
9      ....:                   'F': [datetime.datetime(2020, i, 1) for i in
10      range(1, 13)] +
11      ....:                   [datetime.datetime(2020, i, 15) for i in
12      range(1, 13)]})
13
14  In [72]: df4
15  Out[72]:
16
17      A  B  C      D      E      F
18  0  one  A  foo  0.158188 -0.194267 2020-01-01
19  1  one  B  foo  1.997928  0.087429 2020-02-01
20  2  two  C  foo -2.187545 -0.847917 2020-03-01
21  3  three A  bar  0.448641  1.804847 2020-04-01
22  4  one  B  bar  0.753182  0.349110 2020-05-01
23  5  one  C  bar -0.124078 -1.079033 2020-06-01
24  6  two  A  foo -0.260935  0.977308 2020-07-01
25  7  three B  foo  1.622864  1.303867 2020-08-01
26  8  one  C  foo -0.649257 -0.422324 2020-09-01
27  9  one  A  bar  0.365891 -1.304234 2020-10-01
28 10  two  B  bar  2.759241  0.412113 2020-11-01
29 11  three C  bar  1.247950 -1.616772 2020-12-01
30 12  one  A  foo  0.630572  0.538397 2020-01-15
31 13  one  B  foo -0.404018 -0.239144 2020-02-15
32 14  two  C  foo -0.012304 -0.303686 2020-03-15
33 15  three A  bar -0.510348 -0.513858 2020-04-15
34 16  one  B  bar  1.291382  0.492975 2020-05-15
35 17  one  C  bar  1.272235  0.131740 2020-06-15
36 18  two  A  foo -0.101000 -0.700864 2020-07-15
37 19  three B  foo -2.651767  0.554233 2020-08-15
38 20  one  C  foo  0.714918 -0.489591 2020-09-15
39 21  one  A  bar -0.189745 -0.781274 2020-10-15
40 22  two  B  bar  0.476801 -0.178456 2020-11-15
41 23  three C  bar  1.134483 -0.933998 2020-12-15

```

```

39 In [73]: pd.pivot_table(df4, values="D", index=["A", "B"], columns=["C"])
40 Out[73]:
41 C          bar          foo
42 A    B
43 one  A  0.088073  0.394380
44       B  1.022282  0.796955
45       C  0.574079  0.032831
46 three A -0.030854         NaN
47       B         NaN -0.514451
48       C  1.191217         NaN
49 two   A         NaN -0.180968
50       B  1.618021         NaN
51       C         NaN -1.099925

```

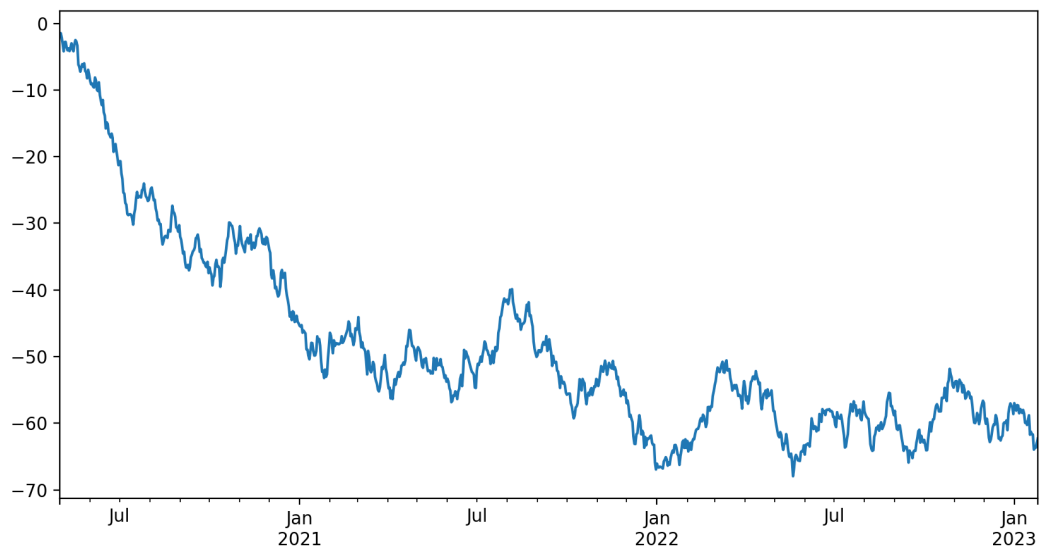
函数 `pivot_table()` 中, 默认 `aggfunc='mean'` 计算均值

Time Series & Graph & File

```

1 # Time Series
2 In [74]: from pylab import *
3
4 In [75]: t_exam = pd.date_range("20200501", periods=10, freq="S")
5
6 In [76]: t_exam
7 Out[76]:
8 DatetimeIndex(['2020-05-01 00:00:00', '2020-05-01 00:00:01',
9                '2020-05-01 00:00:02', '2020-05-01 00:00:03',
10               '2020-05-01 00:00:04', '2020-05-01 00:00:05',
11               '2020-05-01 00:00:06', '2020-05-01 00:00:07',
12               '2020-05-01 00:00:08', '2020-05-01 00:00:09'],
13              dtype='datetime64[ns]', freq='S')
14
15 # Graph
16 In [77]: ts = pd.Series(np.random.randn(1000),
17                        ..., index=pd.date_range("20200501", periods=1000))
18
19 In [78]: ts = ts.cumsum()
20
21 In [79]: ts.plot()
22 Out[79]: <matplotlib.axes._subplots.AxesSubplot at 0x11b166e90>
23
24 In [80]: show()

```



Files

```

1 In [81]: df4.to_csv("datas/Test.csv", index=0)
2
3 In [82]: df4.to_excel("datas/Test.xlsx")
4
5 In [83]: !tree datas
6 datas
7 |— Test.csv
8 |— Test.xlsx
9
10 0 directories, 2 files
11
12 In [84]: df4_csv = pd.read_csv("datas/Test.csv")
13
14 In [85]: df4_csv
15 Out[85]:
16      A  B  C      D      E      F
17 0  one  A  foo  0.158188 -0.194267  2020-01-01
18 1  one  B  foo  1.997928  0.087429  2020-02-01
19 2  two  C  foo -2.187545 -0.847917  2020-03-01
20 ... ..
21 21 one  A  bar -0.189745 -0.781274  2020-10-15
22 22 two  B  bar  0.476801 -0.178456  2020-11-15
23 23 three C  bar  1.134483 -0.933998  2020-12-15
24
25 In [86]: df4_excel = pd.read_excel("datas/Test.xlsx")
26
27 In [87]: df4_excel
28 Out[87]:
29      Unnamed: 0      A  B  C      D      E      F
30 0              0  one  A  foo  0.158188 -0.194267  2020-01-01
31 1              1  one  B  foo  1.997928  0.087429  2020-02-01

```

32	2	2	two	C	foo	-2.187545	-0.847917	2020-03-01
33						
34	21	21	one	A	bar	-0.189745	-0.781274	2020-10-15
35	22	22	two	B	bar	0.476801	-0.178456	2020-11-15
36	23	23	three	C	bar	1.134483	-0.933998	2020-12-15

agC10r