

Pandas Learning

库的导入

```
1 # 导入numpy库并重命名为np  
2 import numpy as np  
3 # 导入pandas 库并重命名为pd  
4 import pandas as pd
```

数据的导入

```
1 pd.read_csv(filename)    # 从csv导入  
2 pd.read_table(filename) # 导入有分隔符的文本（如TSV）中的数据  
3 pd.read_excel(filename) # 从excel导入  
4 pd.read_sql(query, connection_object) # 导入SQL数据表/数据库  
    中的数据  
5 pd.read_json(json_string) # 导入JSON格式的字符，URL地址或者文  
    件中的数据  
6 pd.read_html(url) # 导入经过解析的URL地址中包含的数据框  
    (DataFrame) 数据  
7 pd.read_clipboard() # 导入系统粘贴板里面的数据  
8 pd.DataFrame(dict) # 导入Python字典 (dict) 里面的数据，其中  
    key是数据框的表头，value是数据框的内容。
```

数据的导出

```
1 df.to_csv(filename) # 将数据框 (DataFrame)中的数据导入csv格式  
的文件中  
2 df.to_excel(filename) # 将数据框 (DataFrame)中的数据导入Excel  
格式的文件中  
3 df.to_sql(table_name,connection_object) # 将数据框  
(DataFrame)中的数据导入SQL数据表/数据库中  
4 df.to_json(filename) # 将数据框 (DataFrame)中的数据导入JSON格  
式的文件中
```

创建测试对象

```
1 pd.DataFrame(np.random.rand(5, 10)) # 创建一个5列10行的由随机  
浮点数组成的数据框 DataFrame
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>0</th>  
4   <th>1</th>  
5   <th>2</th>  
6   <th>3</th>  
7   <th>4</th>  
8   <th>5</th>  
9   <th>6</th>  
10  <th>7</th>  
11  <th>8</th>  
12  <th>9</th>  
13 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.016860</td>
4   <td>0.855994</td>
5   <td>0.992872</td>
6   <td>0.652278</td>
7   <td>0.517510</td>
8   <td>0.742986</td>
9   <td>0.452981</td>
10  <td>0.568701</td>
11  <td>0.795436</td>
12  <td>0.622609</td>
13 </tr>
14 <tr>
15   <th>1</th>
16   <td>0.476801</td>
17   <td>0.190823</td>
18   <td>0.450436</td>
19   <td>0.912401</td>
20   <td>0.335651</td>
21   <td>0.197766</td>
22   <td>0.042523</td>
23   <td>0.580323</td>
24   <td>0.498982</td>
25   <td>0.473128</td>
26 </tr>
27 <tr>
28   <th>2</th>
29   <td>0.029820</td>
30   <td>0.886500</td>
31   <td>0.902864</td>
32   <td>0.465084</td>
33   <td>0.380933</td>
34   <td>0.033583</td>
```

```
35 <td>0.928827</td>
36 <td>0.501687</td>
37 <td>0.857512</td>
38 <td>0.671840</td>
39 </tr>
40 <tr>
41 <th>3</th>
42 <td>0.897254</td>
43 <td>0.413717</td>
44 <td>0.991061</td>
45 <td>0.393033</td>
46 <td>0.388630</td>
47 <td>0.661025</td>
48 <td>0.635417</td>
49 <td>0.695609</td>
50 <td>0.305378</td>
51 <td>0.147508</td>
52 </tr>
53 <tr>
54 <th>4</th>
55 <td>0.573882</td>
56 <td>0.786888</td>
57 <td>0.177782</td>
58 <td>0.864474</td>
59 <td>0.594416</td>
60 <td>0.765678</td>
61 <td>0.217279</td>
62 <td>0.446570</td>
63 <td>0.930604</td>
64 <td>0.686823</td>
65 </tr>
```

```
1 pd.Series(my_list) # 从一个可迭代的对象 my_list 中创建一个数据  
组
```

```
1 my_list = ['abc', 123, 'HelloWorld', 5.7]  
2 pd.Series(my_list)
```

```
1 0          abc  
2 1          123  
3 2      HelloWorld  
4 3          5.7  
5 dtype: object
```

```
1 df = pd.DataFrame(np.random.rand(10, 5))  
2 df.index = pd.date_range('2017/1/1', periods=df.shape[0])  
3 df
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>0</th>  
4   <th>1</th>  
5   <th>2</th>  
6   <th>3</th>  
7   <th>4</th>  
8 </tr>
```

```
1 <tr>
2   <th>2017-01-01</th>
3   <td>0.011762</td>
4   <td>0.634116</td>
5   <td>0.045220</td>
6   <td>0.452117</td>
7   <td>0.879969</td>
8 </tr>
9 <tr>
10  <th>2017-01-02</th>
11  <td>0.802262</td>
12  <td>0.661908</td>
13  <td>0.214822</td>
14  <td>0.444259</td>
15  <td>0.200370</td>
16 </tr>
17 <tr>
18  <th>2017-01-03</th>
19  <td>0.301050</td>
20  <td>0.004534</td>
21  <td>0.881042</td>
22  <td>0.825632</td>
23  <td>0.331118</td>
24 </tr>
25 <tr>
26  <th>2017-01-04</th>
27  <td>0.095324</td>
28  <td>0.916430</td>
29  <td>0.177795</td>
30  <td>0.191502</td>
31  <td>0.546973</td>
32 </tr>
33 <tr>
34  <th>2017-01-05</th>
```

```
35 <td>0.482868</td>
36 <td>0.953719</td>
37 <td>0.615461</td>
38 <td>0.868984</td>
39 <td>0.639286</td>
40 </tr>
41 <tr>
42 <th>2017-01-06</th>
43 <td>0.958404</td>
44 <td>0.155357</td>
45 <td>0.293012</td>
46 <td>0.115218</td>
47 <td>0.177846</td>
48 </tr>
49 <tr>
50 <th>2017-01-07</th>
51 <td>0.915488</td>
52 <td>0.486922</td>
53 <td>0.440474</td>
54 <td>0.584764</td>
55 <td>0.271243</td>
56 </tr>
57 <tr>
58 <th>2017-01-08</th>
59 <td>0.480413</td>
60 <td>0.600622</td>
61 <td>0.325212</td>
62 <td>0.532259</td>
63 <td>0.687718</td>
64 </tr>
65 <tr>
66 <th>2017-01-09</th>
67 <td>0.859887</td>
68 <td>0.236677</td>
```

```
69 <td>0.635073</td>
70 <td>0.811840</td>
71 <td>0.497289</td>
72 </tr>
73 <tr>
74 <th>2017-01-10</th>
75 <td>0.024623</td>
76 <td>0.635122</td>
77 <td>0.346393</td>
78 <td>0.860260</td>
79 <td>0.325502</td>
80 </tr>
```

数据的查看

```
1 df.head(n) # 查看前n行的数据
2 df.tail(n) # 查看后n行的数据
```

```
1 df = pd.DataFrame(np.random.rand(10, 5))
2 df.head(3)
```

```
1 <tr style="text-align: right;">>
2   <th></th>
3   <th>0</th>
4   <th>1</th>
5   <th>2</th>
6   <th>3</th>
7   <th>4</th>
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.996381</td>
4   <td>0.440502</td>
5   <td>0.583701</td>
6   <td>0.120444</td>
7   <td>0.241775</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.126877</td>
12  <td>0.646841</td>
13  <td>0.740163</td>
14  <td>0.764182</td>
15  <td>0.810129</td>
16 </tr>
17 <tr>
18  <th>2</th>
19  <td>0.254386</td>
20  <td>0.451341</td>
21  <td>0.288513</td>
22  <td>0.515995</td>
23  <td>0.146529</td>
24 </tr>
```

```
1 df = pd.DataFrame(np.random.rand(10, 5))  
2 df.tail(3)
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>0</th>  
4   <th>1</th>  
5   <th>2</th>  
6   <th>3</th>  
7   <th>4</th>  
8 </tr>
```

```
1 <tr>
2   <th>7</th>
3   <td>0.466316</td>
4   <td>0.747013</td>
5   <td>0.568442</td>
6   <td>0.562552</td>
7   <td>0.949529</td>
8 </tr>
9 <tr>
10  <th>8</th>
11  <td>0.243633</td>
12  <td>0.605133</td>
13  <td>0.114011</td>
14  <td>0.898604</td>
15  <td>0.024648</td>
16 </tr>
17 <tr>
18  <th>9</th>
19  <td>0.155605</td>
20  <td>0.799580</td>
21  <td>0.160883</td>
22  <td>0.986743</td>
23  <td>0.446114</td>
24 </tr>
```

查看数据的形状

```
1 | df.shape # 查看数据的形状（行和宽）
```

查看数据的相关信息

```
1 | df.info() # 查看数据的索引、数据类型及内存信息
```

```
1 | <class 'pandas.core.frame.DataFrame'>
2 | RangeIndex: 10 entries, 0 to 9
3 | Data columns (total 5 columns):
4 | 0    10 non-null float64
5 | 1    10 non-null float64
6 | 2    10 non-null float64
7 | 3    10 non-null float64
8 | 4    10 non-null float64
9 | dtypes: float64(5)
10 | memory usage: 480.0 bytes
```

```
1 | df.describe() # 对于数据类型为数值型的列，查询其描述性统计的内容
```

```
1 | <tr style="text-align: right;">
2 |   <th></th>
3 |   <th>0</th>
4 |   <th>1</th>
5 |   <th>2</th>
6 |   <th>3</th>
7 |   <th>4</th>
8 | </tr>
```

```
1 <tr>
2   <th>count</th>
3   <td>10.000000</td>
4   <td>10.000000</td>
5   <td>10.000000</td>
6   <td>10.000000</td>
7   <td>10.000000</td>
8 </tr>
9 <tr>
10  <th>mean</th>
11  <td>0.432296</td>
12  <td>0.601318</td>
13  <td>0.444123</td>
14  <td>0.543053</td>
15  <td>0.413944</td>
16 </tr>
17 <tr>
18  <th>std</th>
19  <td>0.268632</td>
20  <td>0.216483</td>
21  <td>0.298076</td>
22  <td>0.284759</td>
23  <td>0.345554</td>
24 </tr>
25 <tr>
26  <th>min</th>
27  <td>0.155605</td>
28  <td>0.252712</td>
29  <td>0.114011</td>
30  <td>0.126172</td>
31  <td>0.024648</td>
32 </tr>
33 <tr>
34  <th>25%</th>
```

```
35 <td>0.245740</td>
36 <td>0.461876</td>
37 <td>0.191090</td>
38 <td>0.295224</td>
39 <td>0.094655</td>
40 </tr>
41 <tr>
42 <th>50%</th>
43 <td>0.353927</td>
44 <td>0.609304</td>
45 <td>0.401004</td>
46 <td>0.565072</td>
47 <td>0.389053</td>
48 </tr>
49 <tr>
50 <th>75%</th>
51 <td>0.559964</td>
52 <td>0.762463</td>
53 <td>0.556721</td>
54 <td>0.693000</td>
55 <td>0.668621</td>
56 </tr>
57 <tr>
58 <th>max</th>
59 <td>0.985457</td>
60 <td>0.932679</td>
61 <td>0.990827</td>
62 <td>0.986743</td>
63 <td>0.949529</td>
64 </tr>
```

统计次数

```
1 | s.value_counts(dropna=False) # 查询每个独特数据值出现次数统计
```

```
1 | s = pd.Series([1,2,3,3,4,np.nan,5,5,5,6,7])
2 | s.value_counts(dropna=False)
```

```
1 |      5.0    3
2 |      3.0    2
3 |      7.0    1
4 |      6.0    1
5 |      NaN    1
6 |      4.0    1
7 |      2.0    1
8 |      1.0    1
9 |      dtype: int64
```

```
1 | s = s.apply(lambda x: x+1)
2 | s
```

```
1 |      0    3.0
2 |      1    4.0
3 |      2    5.0
4 |      3    5.0
5 |      4    6.0
6 |      5    NaN
7 |      6    7.0
8 |      7    7.0
9 |      8    7.0
10 |     9    8.0
11 |    10    9.0
12 |      dtype: float64
```

```
1 df.apply(pd.Series.value_counts) # 查询数据框 (Data Frame) 中  
每个列的独特数据值出现次数统计
```

数据选取

```
1 df[col] # 以数组 Series 的形式返回选取的列
```

```
1 df = pd.DataFrame(np.random.rand(5, 5),  
columns=list('ABCDE'))  
2 df['C']
```

```
1 0    0.452717  
2 1    0.407755  
3 2    0.549391  
4 3    0.759433  
5 4    0.153871  
6 Name: C, dtype: float64
```

```
1 df[[col1, col2]] # 选择多列
```

```
1 df = pd.DataFrame(np.random.rand(5, 5),  
columns=list('ABCDE'))  
2 df[['C', 'D']]
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>C</th>  
4   <th>D</th>  
5 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.431885</td>
4   <td>0.304796</td>
5 </tr>
6 <tr>
7   <th>1</th>
8   <td>0.028960</td>
9   <td>0.187738</td>
10 </tr>
11 <tr>
12   <th>2</th>
13   <td>0.176520</td>
14   <td>0.102980</td>
15 </tr>
16 <tr>
17   <th>3</th>
18   <td>0.370277</td>
19   <td>0.098031</td>
20 </tr>
21 <tr>
22   <th>4</th>
23   <td>0.247122</td>
24   <td>0.345735</td>
25 </tr>
```

```
1 | s.iloc[0] # 按位置选取
```

```
1 | s = pd.Series(np.array(['I', 'Love', 'China']))
2 | s.iloc[0]
```

```
1 | 'I'
```

```
1 | s.loc['index_one'] # 按索引选取
```

```
1 | s = pd.Series(np.array(['I', 'Love', 'China']))
2 | s.loc[0]
```

```
1 | 'I'
```

```
1 | df.DataFrame[n, :] #选取第n行
```

```
1 | df = pd.DataFrame(np.array([[['I', 'Love', 'China'], ['I',
2 | 'Love', 'Data']])))
2 | df.iloc[1, :]
```

```
1 | 0      I
2 | 1      Love
3 | 2      Data
4 | Name: 1, dtype: object
```

```
1 | df.iloc[0, 0] # 选取第一个元素
```

```
1 | df = pd.DataFrame(np.random.rand(5, 5))
2 | df
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>0</th>
4   <th>1</th>
5   <th>2</th>
6   <th>3</th>
7   <th>4</th>
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.786709</td>
4   <td>0.405902</td>
5   <td>0.151383</td>
6   <td>0.384778</td>
7   <td>0.871664</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.491006</td>
12  <td>0.774710</td>
13  <td>0.388011</td>
14  <td>0.758102</td>
15  <td>0.762115</td>
16 </tr>
17 <tr>
18  <th>2</th>
19  <td>0.085647</td>
20  <td>0.543243</td>
21  <td>0.582565</td>
22  <td>0.664243</td>
23  <td>0.379896</td>
24 </tr>
25 <tr>
26  <th>3</th>
27  <td>0.806211</td>
28  <td>0.794284</td>
29  <td>0.968755</td>
30  <td>0.883923</td>
31  <td>0.354820</td>
32 </tr>
33 <tr>
34  <th>4</th>
```

```
35 <td>0.463902</td>
36 <td>0.481756</td>
37 <td>0.131181</td>
38 <td>0.590878</td>
39 <td>0.801769</td>
40 </tr>
```

```
1 df.iloc[0, 0]
```

```
1 0.78670886755075187
```

数据的清洗

```
1 df.columns = ['a', 'b'] # 对列名重新命名
```

```
1 df = pd.DataFrame({'A':np.array([1,np.nan,2,3,6,np.nan]),
2
3     'B':np.array([np.nan,4,np.nan,5,9,np.nan]),
4         'C':'foo'})
```

```
1 <tr style="text-align: right;">
2     <th></th>
3     <th>A</th>
4     <th>B</th>
5     <th>C</th>
6 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>1.0</td>
4   <td>NaN</td>
5   <td>foo</td>
6 </tr>
7 <tr>
8   <th>1</th>
9   <td>NaN</td>
10  <td>4.0</td>
11  <td>foo</td>
12 </tr>
13 <tr>
14   <th>2</th>
15   <td>2.0</td>
16   <td>NaN</td>
17   <td>foo</td>
18 </tr>
19 <tr>
20   <th>3</th>
21   <td>3.0</td>
22   <td>5.0</td>
23   <td>foo</td>
24 </tr>
25 <tr>
26   <th>4</th>
27   <td>6.0</td>
28   <td>9.0</td>
29   <td>foo</td>
30 </tr>
31 <tr>
32   <th>5</th>
33   <td>NaN</td>
34   <td>NaN</td>
```

```
35 |     <td>foo</td>
36 | </tr>
```

```
1 | df.columns = ['a', 'b', 'c']
2 | df
```

```
1 | <tr style="text-align: right;">
2 |     <th></th>
3 |     <th>a</th>
4 |     <th>b</th>
5 |     <th>c</th>
6 | </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>1.0</td>
4   <td>NaN</td>
5   <td>foo</td>
6 </tr>
7 <tr>
8   <th>1</th>
9   <td>NaN</td>
10  <td>4.0</td>
11  <td>foo</td>
12 </tr>
13 <tr>
14   <th>2</th>
15   <td>2.0</td>
16   <td>NaN</td>
17   <td>foo</td>
18 </tr>
19 <tr>
20   <th>3</th>
21   <td>3.0</td>
22   <td>5.0</td>
23   <td>foo</td>
24 </tr>
25 <tr>
26   <th>4</th>
27   <td>6.0</td>
28   <td>9.0</td>
29   <td>foo</td>
30 </tr>
31 <tr>
32   <th>5</th>
33   <td>NaN</td>
34   <td>NaN</td>
```

```
35 <td>foo</td>
36 </tr>
```

```
1 pd.isnull() # 检查数据中出现空值的情况，返回一个布尔型的列
2 pd.notnull() #相对应isnull 返回不是空值的情况
```

```
1 df = pd.DataFrame({'A':np.array([1,np.nan,2,3,6,np.nan]),
2
3     'B':np.array([np.nan,4,np.nan,5,9,np.nan]),
4         'C':'foo'})
```

```
1 <tr style="text-align: right;">
2     <th></th>
3     <th>A</th>
4     <th>B</th>
5     <th>C</th>
6 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>False</td>
4   <td>True</td>
5   <td>False</td>
6 </tr>
7 <tr>
8   <th>1</th>
9   <td>True</td>
10  <td>False</td>
11  <td>False</td>
12 </tr>
13 <tr>
14   <th>2</th>
15   <td>False</td>
16   <td>True</td>
17   <td>False</td>
18 </tr>
19 <tr>
20   <th>3</th>
21   <td>False</td>
22   <td>False</td>
23   <td>False</td>
24 </tr>
25 <tr>
26   <th>4</th>
27   <td>False</td>
28   <td>False</td>
29   <td>False</td>
30 </tr>
31 <tr>
32   <th>5</th>
33   <td>True</td>
34   <td>True</td>
```

```
35 |     <td>False</td>
36 | </tr>
```

```
1 | df.isnull().sum() # 对每一列的空值进行统计
```

```
1 | A    2
2 | B    3
3 | C    0
4 | dtype: int64
```

```
1 | df.dropna(axis = 0, thresh=n) # 删除包含缺失值的行 axis = 1时
|   删除列 # thresh = n 移除空值超过(包括等于)n的行
```

```
1 | df
```

```
1 | <tr style="text-align: right;">
2 |     <th></th>
3 |     <th>A</th>
4 |     <th>B</th>
5 |     <th>C</th>
6 | </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>1.0</td>
4   <td>NaN</td>
5   <td>foo</td>
6 </tr>
7 <tr>
8   <th>1</th>
9   <td>NaN</td>
10  <td>4.0</td>
11  <td>foo</td>
12 </tr>
13 <tr>
14   <th>2</th>
15   <td>2.0</td>
16   <td>NaN</td>
17   <td>foo</td>
18 </tr>
19 <tr>
20   <th>3</th>
21   <td>3.0</td>
22   <td>5.0</td>
23   <td>foo</td>
24 </tr>
25 <tr>
26   <th>4</th>
27   <td>6.0</td>
28   <td>9.0</td>
29   <td>foo</td>
30 </tr>
31 <tr>
32   <th>5</th>
33   <td>NaN</td>
34   <td>NaN</td>
```

```
35 |     <td>foo</td>
36 | </tr>
```

```
1 | df.dropna(axis = 0)
```

```
1 | <tr style="text-align: right;">
2 |     <th></th>
3 |     <th>A</th>
4 |     <th>B</th>
5 |     <th>C</th>
6 | </tr>
```

```
1 | <tr>
2 |     <th>3</th>
3 |     <td>3.0</td>
4 |     <td>5.0</td>
5 |     <td>foo</td>
6 | </tr>
7 | <tr>
8 |     <th>4</th>
9 |     <td>6.0</td>
10 |    <td>9.0</td>
11 |    <td>foo</td>
12 | </tr>
```

```
1 | df.dropna(axis = 1)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>C</th>
4 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>foo</td>
4 </tr>
5 <tr>
6   <th>1</th>
7   <td>foo</td>
8 </tr>
9 <tr>
10  <th>2</th>
11  <td>foo</td>
12 </tr>
13 <tr>
14  <th>3</th>
15  <td>foo</td>
16 </tr>
17 <tr>
18  <th>4</th>
19  <td>foo</td>
20 </tr>
21 <tr>
22  <th>5</th>
23  <td>foo</td>
24 </tr>
```

```
1 df.dropna(axis = 0, thresh = 2)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>1.0</td>
4   <td>NaN</td>
5   <td>foo</td>
6 </tr>
7 <tr>
8   <th>1</th>
9   <td>NaN</td>
10  <td>4.0</td>
11  <td>foo</td>
12 </tr>
13 <tr>
14   <th>2</th>
15   <td>2.0</td>
16   <td>NaN</td>
17   <td>foo</td>
18 </tr>
19 <tr>
20   <th>3</th>
21   <td>3.0</td>
22   <td>5.0</td>
23   <td>foo</td>
24 </tr>
25 <tr>
26   <th>4</th>
27   <td>6.0</td>
28   <td>9.0</td>
29   <td>foo</td>
30 </tr>
```

```
1 df.fillna(df.mean()) # 用平均值来填充空值
```

```
1 s = pd.Series([1,3,5,np.nan,7,9,9])
2 s.fillna(s.mean())
```

```
1 0    1.000000
2 1    3.000000
3 2    5.000000
4 3    5.666667
5 4    7.000000
6 5    9.000000
7 6    9.000000
8 dtype: float64
```

```
1 s.astype(type) # 转换列的类型
```

```
1 s = pd.Series([1,3,5,np.nan,7,9,9])
2 s.fillna(s.mean()).astype(int)
```

```
1 0    1
2 1    3
3 2    5
4 3    5
5 4    7
6 5    9
7 6    9
8 dtype: int32
```

```
1 s.replace(1, 'one') # 将Series中的1替换为one
```

```
1 s = pd.Series([1,3,5,np.nan,7,9,9])
2 s.replace(1,'one')
```

```
1 0      one
2 1      3
3 2      5
4 3      NaN
5 4      7
6 5      9
7 6      9
8 dtype: object
```

```
1 s.replace([1,3],['one','three']) # 将数组(Series)中所有的1替
换为'one'，所有的3替换为'three'
```

```
1 s = pd.Series([1,3,5,np.nan,7,9,9])
2 s.replace([1,3],['one','three'])
```

```
1 0      one
2 1      three
3 2      5
4 3      NaN
5 4      7
6 5      9
7 6      9
8 dtype: object
```

```
1 df.rename(columns=lambda x: x + 2) # 将全体列重命名
```

```
1 df = pd.DataFrame(np.random.rand(4,4))
2 df
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>0</th>
4   <th>1</th>
5   <th>2</th>
6   <th>3</th>
7 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.669102</td>
4   <td>0.548996</td>
5   <td>0.512802</td>
6   <td>0.449220</td>
7 </tr>
8 <tr>
9   <th>1</th>
10  <td>0.108840</td>
11  <td>0.974720</td>
12  <td>0.665050</td>
13  <td>0.271009</td>
14 </tr>
15 <tr>
16   <th>2</th>
17   <td>0.146804</td>
18   <td>0.060744</td>
19   <td>0.637770</td>
20   <td>0.383380</td>
21 </tr>
22 <tr>
23   <th>3</th>
24   <td>0.108163</td>
25   <td>0.893999</td>
26   <td>0.216907</td>
27   <td>0.730504</td>
28 </tr>
```

```
1 | df.rename(columns=lambda x: x+ 2)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>2</th>
4   <th>3</th>
5   <th>4</th>
6   <th>5</th>
7 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.669102</td>
4   <td>0.548996</td>
5   <td>0.512802</td>
6   <td>0.449220</td>
7 </tr>
8 <tr>
9   <th>1</th>
10  <td>0.108840</td>
11  <td>0.974720</td>
12  <td>0.665050</td>
13  <td>0.271009</td>
14 </tr>
15 <tr>
16   <th>2</th>
17   <td>0.146804</td>
18   <td>0.060744</td>
19   <td>0.637770</td>
20   <td>0.383380</td>
21 </tr>
22 <tr>
23   <th>3</th>
24   <td>0.108163</td>
25   <td>0.893999</td>
26   <td>0.216907</td>
27   <td>0.730504</td>
28 </tr>
```

```
1 df.rename(columns={'old_name': 'new_ name'}) # 将选择的列重命名
```

```
1 df =  
pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))  
2 df.rename(columns={'A':'New A', 'B':'New B'})
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>New A</th>  
4   <th>New B</th>  
5   <th>C</th>  
6   <th>D</th>  
7   <th>E</th>  
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.680941</td>
4   <td>0.766561</td>
5   <td>0.486226</td>
6   <td>0.301537</td>
7   <td>0.289970</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.917036</td>
12  <td>0.100054</td>
13  <td>0.464342</td>
14  <td>0.181454</td>
15  <td>0.933591</td>
16 </tr>
17 <tr>
18  <th>2</th>
19  <td>0.253549</td>
20  <td>0.766181</td>
21  <td>0.085607</td>
22  <td>0.969627</td>
23  <td>0.630674</td>
24 </tr>
25 <tr>
26  <th>3</th>
27  <td>0.377840</td>
28  <td>0.909920</td>
29  <td>0.214338</td>
30  <td>0.011844</td>
31  <td>0.392257</td>
32 </tr>
33 <tr>
34  <th>4</th>
```

```
35 <td>0.608564</td>
36 <td>0.587614</td>
37 <td>0.039867</td>
38 <td>0.630492</td>
39 <td>0.402101</td>
40 </tr>
41 <tr>
42 <th>5</th>
43 <td>0.361074</td>
44 <td>0.937618</td>
45 <td>0.787055</td>
46 <td>0.054157</td>
47 <td>0.300325</td>
48 </tr>
49 <tr>
50 <th>6</th>
51 <td>0.605472</td>
52 <td>0.608429</td>
53 <td>0.052152</td>
54 <td>0.669343</td>
55 <td>0.745648</td>
56 </tr>
57 <tr>
58 <th>7</th>
59 <td>0.660738</td>
60 <td>0.158713</td>
61 <td>0.352756</td>
62 <td>0.028325</td>
63 <td>0.195899</td>
64 </tr>
65 <tr>
66 <th>8</th>
67 <td>0.855695</td>
68 <td>0.578177</td>
```

```
69 <td>0.447043</td>
70 <td>0.093923</td>
71 <td>0.316234</td>
72 </tr>
73 <tr>
74 <th>9</th>
75 <td>0.337392</td>
76 <td>0.645260</td>
77 <td>0.140221</td>
78 <td>0.616652</td>
79 <td>0.727144</td>
80 </tr>
```

```
1 df.set_index('column_one') # 改变索引
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 df.set_index('B')
```

```
1 <tr style="text-align: right;">>
2   <th></th>
3   <th>A</th>
4   <th>C</th>
5   <th>D</th>
6   <th>E</th>
7 </tr>
8 <tr>
9   <th>B</th>
10  <th></th>
11  <th></th>
12  <th></th>
13  <th></th>
14 </tr>
```

```
1 <tr>
2   <th>0.824403</th>
3   <td>0.013092</td>
4   <td>0.055626</td>
5   <td>0.895268</td>
6   <td>0.837350</td>
7 </tr>
8 <tr>
9   <th>0.310559</th>
10  <td>0.689064</td>
11  <td>0.541375</td>
12  <td>0.275461</td>
13  <td>0.808554</td>
14 </tr>
15 <tr>
16   <th>0.533495</th>
17   <td>0.072835</td>
18   <td>0.563758</td>
19   <td>0.695029</td>
20   <td>0.524957</td>
21 </tr>
22 <tr>
23   <th>0.541211</th>
24   <td>0.820817</td>
25   <td>0.591130</td>
26   <td>0.268978</td>
27   <td>0.546996</td>
28 </tr>
29 <tr>
30   <th>0.286587</th>
31   <td>0.936692</td>
32   <td>0.343227</td>
33   <td>0.383610</td>
34   <td>0.811302</td>
```

```
35 </tr>
36 <tr>
37   <th>0.878391</th>
38   <td>0.938883</td>
39   <td>0.636148</td>
40   <td>0.776493</td>
41   <td>0.025840</td>
42 </tr>
43 <tr>
44   <th>0.156482</th>
45   <td>0.918591</td>
46   <td>0.030869</td>
47   <td>0.235020</td>
48   <td>0.096212</td>
49 </tr>
50 <tr>
51   <th>0.857049</th>
52   <td>0.613991</td>
53   <td>0.810541</td>
54   <td>0.917927</td>
55   <td>0.921329</td>
56 </tr>
57 <tr>
58   <th>0.713271</th>
59   <td>0.949683</td>
60   <td>0.811386</td>
61   <td>0.920452</td>
62   <td>0.213173</td>
63 </tr>
64 <tr>
65   <th>0.686945</th>
66   <td>0.522276</td>
67   <td>0.881299</td>
68   <td>0.936260</td>
```

```
69 |     <td>0.030993</td>
70 | </tr>
```

```
1 | df.rename(index = lambda x: x+ 1) # 改变全体索引
```

```
1 | df = pd.DataFrame(np.random.rand(10,5))
2 | df.rename(index = lambda x: x+ 1)
```

```
1 | <tr style="text-align: right;">
2 |     <th></th>
3 |     <th>0</th>
4 |     <th>1</th>
5 |     <th>2</th>
6 |     <th>3</th>
7 |     <th>4</th>
8 | </tr>
```

```
1 <tr>
2   <th>1</th>
3   <td>0.382337</td>
4   <td>0.185501</td>
5   <td>0.457958</td>
6   <td>0.009713</td>
7   <td>0.628963</td>
8 </tr>
9 <tr>
10  <th>2</th>
11  <td>0.024175</td>
12  <td>0.223274</td>
13  <td>0.698171</td>
14  <td>0.071715</td>
15  <td>0.063272</td>
16 </tr>
17 <tr>
18  <th>3</th>
19  <td>0.913995</td>
20  <td>0.713092</td>
21  <td>0.269621</td>
22  <td>0.575365</td>
23  <td>0.805266</td>
24 </tr>
25 <tr>
26  <th>4</th>
27  <td>0.612708</td>
28  <td>0.220953</td>
29  <td>0.090858</td>
30  <td>0.425472</td>
31  <td>0.018996</td>
32 </tr>
33 <tr>
34  <th>5</th>
```

```
35 <td>0.045363</td>
36 <td>0.153343</td>
37 <td>0.730828</td>
38 <td>0.323554</td>
39 <td>0.364821</td>
40 </tr>
41 <tr>
42 <th>6</th>
43 <td>0.462096</td>
44 <td>0.614072</td>
45 <td>0.993130</td>
46 <td>0.988894</td>
47 <td>0.788648</td>
48 </tr>
49 <tr>
50 <th>7</th>
51 <td>0.887381</td>
52 <td>0.802119</td>
53 <td>0.191248</td>
54 <td>0.980064</td>
55 <td>0.628450</td>
56 </tr>
57 <tr>
58 <th>8</th>
59 <td>0.138270</td>
60 <td>0.922870</td>
61 <td>0.250827</td>
62 <td>0.297472</td>
63 <td>0.289915</td>
64 </tr>
65 <tr>
66 <th>9</th>
67 <td>0.258687</td>
68 <td>0.807993</td>
```

```
69 <td>0.930009</td>
70 <td>0.811335</td>
71 <td>0.609763</td>
72 </tr>
73 <tr>
74 <th>10</th>
75 <td>0.588020</td>
76 <td>0.392127</td>
77 <td>0.590799</td>
78 <td>0.923180</td>
79 <td>0.722801</td>
80 </tr>
```

数据的过滤(filter),排序(sort)和分组(groupby)

```
1 df[df[col] > 0.5] # 选取数据df中对应行的数值大于0.5的全部列 支持逻辑运算
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
2 df
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7   <th>E</th>
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.051900</td>
4   <td>0.548808</td>
5   <td>0.744936</td>
6   <td>0.848002</td>
7   <td>0.299505</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.979053</td>
12  <td>0.216078</td>
13  <td>0.394286</td>
14  <td>0.520654</td>
15  <td>0.584194</td>
16 </tr>
17 <tr>
18  <th>2</th>
19  <td>0.185679</td>
20  <td>0.453151</td>
21  <td>0.839947</td>
22  <td>0.730177</td>
23  <td>0.392377</td>
24 </tr>
25 <tr>
26  <th>3</th>
27  <td>0.161267</td>
28  <td>0.981833</td>
29  <td>0.890858</td>
30  <td>0.613972</td>
31  <td>0.467528</td>
32 </tr>
33 <tr>
34  <th>4</th>
```

```
35 <td>0.091140</td>
36 <td>0.369805</td>
37 <td>0.600035</td>
38 <td>0.372857</td>
39 <td>0.897063</td>
40 </tr>
41 <tr>
42 <th>5</th>
43 <td>0.612195</td>
44 <td>0.981150</td>
45 <td>0.578304</td>
46 <td>0.220064</td>
47 <td>0.488182</td>
48 </tr>
49 <tr>
50 <th>6</th>
51 <td>0.898736</td>
52 <td>0.626289</td>
53 <td>0.788306</td>
54 <td>0.747086</td>
55 <td>0.386097</td>
56 </tr>
57 <tr>
58 <th>7</th>
59 <td>0.568531</td>
60 <td>0.362593</td>
61 <td>0.644950</td>
62 <td>0.510410</td>
63 <td>0.092556</td>
64 </tr>
65 <tr>
66 <th>8</th>
67 <td>0.872898</td>
68 <td>0.771917</td>
```

```
69 <td>0.853365</td>
70 <td>0.227531</td>
71 <td>0.045184</td>
72 </tr>
73 <tr>
74 <th>9</th>
75 <td>0.898296</td>
76 <td>0.683850</td>
77 <td>0.138142</td>
78 <td>0.956854</td>
79 <td>0.335476</td>
80 </tr>
```

```
1 | df[df['A'] > 0.5]
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7   <th>E</th>
8 </tr>
```

```
1 <tr>
2   <th>1</th>
3   <td>0.979053</td>
4   <td>0.216078</td>
5   <td>0.394286</td>
6   <td>0.520654</td>
7   <td>0.584194</td>
8 </tr>
9 <tr>
10  <th>5</th>
11  <td>0.612195</td>
12  <td>0.981150</td>
13  <td>0.578304</td>
14  <td>0.220064</td>
15  <td>0.488182</td>
16 </tr>
17 <tr>
18  <th>6</th>
19  <td>0.898736</td>
20  <td>0.626289</td>
21  <td>0.788306</td>
22  <td>0.747086</td>
23  <td>0.386097</td>
24 </tr>
25 <tr>
26  <th>7</th>
27  <td>0.568531</td>
28  <td>0.362593</td>
29  <td>0.644950</td>
30  <td>0.510410</td>
31  <td>0.092556</td>
32 </tr>
33 <tr>
34  <th>8</th>
```

```
35 <td>0.872898</td>
36 <td>0.771917</td>
37 <td>0.853365</td>
38 <td>0.227531</td>
39 <td>0.045184</td>
40 </tr>
41 <tr>
42 <th>9</th>
43 <td>0.898296</td>
44 <td>0.683850</td>
45 <td>0.138142</td>
46 <td>0.956854</td>
47 <td>0.335476</td>
48 </tr>
```

```
1 | df[(df['A'] > 0.5) & (df['B'] < 0.7)]
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7   <th>E</th>
8 </tr>
```

```
1 <tr>
2   <th>1</th>
3   <td>0.979053</td>
4   <td>0.216078</td>
5   <td>0.394286</td>
6   <td>0.520654</td>
7   <td>0.584194</td>
8 </tr>
9 <tr>
10  <th>6</th>
11  <td>0.898736</td>
12  <td>0.626289</td>
13  <td>0.788306</td>
14  <td>0.747086</td>
15  <td>0.386097</td>
16 </tr>
17 <tr>
18  <th>7</th>
19  <td>0.568531</td>
20  <td>0.362593</td>
21  <td>0.644950</td>
22  <td>0.510410</td>
23  <td>0.092556</td>
24 </tr>
25 <tr>
26  <th>9</th>
27  <td>0.898296</td>
28  <td>0.683850</td>
29  <td>0.138142</td>
30  <td>0.956854</td>
31  <td>0.335476</td>
32 </tr>
```

排序

```
1 df.sort_values(col, ascending=True) #按照列进行排序 #  
ascending: True 升序 False 降序
```

```
1 df =  
pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))  
2 df.sort_values('A', ascending=True)
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>A</th>  
4   <th>B</th>  
5   <th>C</th>  
6   <th>D</th>  
7   <th>E</th>  
8 </tr>
```

```
1 <tr>
2   <th>3</th>
3   <td>0.015834</td>
4   <td>0.758417</td>
5   <td>0.123415</td>
6   <td>0.802403</td>
7   <td>0.782450</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.068046</td>
12  <td>0.373240</td>
13  <td>0.414358</td>
14  <td>0.105285</td>
15  <td>0.759001</td>
16 </tr>
17 <tr>
18  <th>7</th>
19  <td>0.134238</td>
20  <td>0.104416</td>
21  <td>0.551595</td>
22  <td>0.472277</td>
23  <td>0.015997</td>
24 </tr>
25 <tr>
26  <th>4</th>
27  <td>0.236628</td>
28  <td>0.391852</td>
29  <td>0.390275</td>
30  <td>0.904988</td>
31  <td>0.650108</td>
32 </tr>
33 <tr>
34  <th>8</th>
```

```
35 <td>0.469382</td>
36 <td>0.426359</td>
37 <td>0.137109</td>
38 <td>0.253183</td>
39 <td>0.894667</td>
40 </tr>
41 <tr>
42 <th>2</th>
43 <td>0.508937</td>
44 <td>0.443894</td>
45 <td>0.147076</td>
46 <td>0.149885</td>
47 <td>0.434802</td>
48 </tr>
49 <tr>
50 <th>0</th>
51 <td>0.572640</td>
52 <td>0.369032</td>
53 <td>0.412343</td>
54 <td>0.402019</td>
55 <td>0.445365</td>
56 </tr>
57 <tr>
58 <th>9</th>
59 <td>0.663964</td>
60 <td>0.533604</td>
61 <td>0.217605</td>
62 <td>0.602667</td>
63 <td>0.637232</td>
64 </tr>
65 <tr>
66 <th>6</th>
67 <td>0.765109</td>
68 <td>0.646277</td>
```

```
69 <td>0.885381</td>
70 <td>0.743307</td>
71 <td>0.649711</td>
72 </tr>
73 <tr>
74 <th>5</th>
75 <td>0.962494</td>
76 <td>0.650830</td>
77 <td>0.754514</td>
78 <td>0.578115</td>
79 <td>0.659846</td>
80 </tr>
```

```
1 df.sort_values([col1,col2],ascending=[True,False]) # 按照数据框的列col1升序, col2降序的方式对数据框df做排序
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 df.sort_values(['A','E'],ascending=[True,False])
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7   <th>E</th>
8 </tr>
```

```
1 <tr>
2   <th>1</th>
3   <td>0.083526</td>
4   <td>0.969267</td>
5   <td>0.012550</td>
6   <td>0.672851</td>
7   <td>0.866501</td>
8 </tr>
9 <tr>
10  <th>7</th>
11  <td>0.107839</td>
12  <td>0.383900</td>
13  <td>0.982337</td>
14  <td>0.390914</td>
15  <td>0.308559</td>
16 </tr>
17 <tr>
18  <th>0</th>
19  <td>0.148742</td>
20  <td>0.306888</td>
21  <td>0.853949</td>
22  <td>0.144650</td>
23  <td>0.414474</td>
24 </tr>
25 <tr>
26  <th>4</th>
27  <td>0.190011</td>
28  <td>0.794060</td>
29  <td>0.514756</td>
30  <td>0.272207</td>
31  <td>0.086894</td>
32 </tr>
33 <tr>
34  <th>6</th>
```

```
35 <td>0.423982</td>
36 <td>0.229873</td>
37 <td>0.992318</td>
38 <td>0.495706</td>
39 <td>0.971735</td>
40 </tr>
41 <tr>
42 <th>9</th>
43 <td>0.532263</td>
44 <td>0.106900</td>
45 <td>0.528114</td>
46 <td>0.456583</td>
47 <td>0.362642</td>
48 </tr>
49 <tr>
50 <th>2</th>
51 <td>0.619678</td>
52 <td>0.800373</td>
53 <td>0.927766</td>
54 <td>0.742667</td>
55 <td>0.645809</td>
56 </tr>
57 <tr>
58 <th>3</th>
59 <td>0.815312</td>
60 <td>0.920682</td>
61 <td>0.833351</td>
62 <td>0.266840</td>
63 <td>0.132698</td>
64 </tr>
65 <tr>
66 <th>5</th>
67 <td>0.842293</td>
68 <td>0.049499</td>
```

```
69   <td>0.780198</td>
70   <td>0.343752</td>
71   <td>0.341800</td>
72 </tr>
73 <tr>
74   <th>8</th>
75   <td>0.932379</td>
76   <td>0.398721</td>
77   <td>0.080358</td>
78   <td>0.200681</td>
79   <td>0.549237</td>
80 </tr>
```

分组

```
1 df.groupby(col) # 按照某列对数据框df做分组 # 常与count进行连
用，统计出各词的个数
```

```
1 df =
2 pd.DataFrame({ 'A':np.array(['foo','foo','foo','foo','bar','b
ar']),
3
4     'B':np.array(['one','one','two','two','three','three']),
5
6     'C':np.array(['small','medium','large','large','small','smal
l']),
7
8     'D':np.array([1,2,2,3,3,5])})
9 df
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>foo</td>
4   <td>one</td>
5   <td>small</td>
6   <td>1</td>
7 </tr>
8 <tr>
9   <th>1</th>
10  <td>foo</td>
11  <td>one</td>
12  <td>medium</td>
13  <td>2</td>
14 </tr>
15 <tr>
16   <th>2</th>
17   <td>foo</td>
18   <td>two</td>
19   <td>large</td>
20   <td>2</td>
21 </tr>
22 <tr>
23   <th>3</th>
24   <td>foo</td>
25   <td>two</td>
26   <td>large</td>
27   <td>3</td>
28 </tr>
29 <tr>
30   <th>4</th>
31   <td>bar</td>
32   <td>three</td>
33   <td>small</td>
34   <td>3</td>
```

```
35 </tr>
36 <tr>
37   <th>5</th>
38   <td>bar</td>
39   <td>three</td>
40   <td>small</td>
41   <td>5</td>
42 </tr>
```

```
1 | df.groupby('B').count()
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>C</th>
5   <th>D</th>
6 </tr>
7 <tr>
8   <th>B</th>
9   <th></th>
10  <th></th>
11  <th></th>
12 </tr>
```

```
1 <tr>
2   <th>one</th>
3   <td>2</td>
4   <td>2</td>
5   <td>2</td>
6 </tr>
7 <tr>
8   <th>three</th>
9   <td>2</td>
10  <td>2</td>
11  <td>2</td>
12 </tr>
13 <tr>
14   <th>two</th>
15   <td>2</td>
16   <td>2</td>
17   <td>2</td>
18 </tr>
```

```
1 | df.groupby([col1,col2]) # 按照列col1和col2对数据框df做分组
```

```
1 df =
2     pd.DataFrame({ 'A':np.array(['foo','foo','foo','foo','bar','bar']),
3
4         'B':np.array(['one','one','two','two','three','three']),
5
6         'C':np.array(['small','medium','large','large','small','small']),
7
8         'D':np.array([1,2,2,3,3,5]))}
9 df.groupby(['A', 'B']).sum()
```

```
1 <tr style="text-align: right;">
2     <th></th>
3     <th></th>
4     <th>D</th>
5 </tr>
6 <tr>
7     <th>A</th>
8     <th>B</th>
9     <th></th>
10 </tr>
```

```
1 <tr>
2   <th>bar</th>
3   <th>three</th>
4   <td>8</td>
5 </tr>
6 <tr>
7   <th rowspan="2" valign="top">foo</th>
8   <th>one</th>
9   <td>3</td>
10 </tr>
11 <tr>
12   <th>two</th>
13   <td>5</td>
14 </tr>
```

```
1 df.groupby('A')[['B']].mean() # 按照列A对数据框df做分组处理
    后，返回对应的B的平均值
```

```
1 df =
2 pd.DataFrame({'A':np.array(['foo','foo','foo','foo','bar','bar']),
3
4     'B':np.array(['one','one','two','two','three','three']),
5
6     'C':np.array(['small','medium','large','large','small','small']),
7     'D':np.array([1,2,2,3,3,5])})
8 df.groupby('A')[['D']].mean()
```

```
1 A
2 bar    4
3 foo    2
4 Name: D, dtype: int32
```

```
1 df.pivot_table(index=col1,values=[col2,col3],aggfunc=mean) #  
做透视表，索引为col1,针对的数值列为col2和col3，分组函数为平均值
```

```
1 df =
2 pd.DataFrame({'A':np.array(['foo','foo','foo','foo','bar','bar']),
3
4     'B':np.array(['one','one','two','two','three','three']),
5
6     'C':np.array(['small','medium','large','large','small','small']),
7     'D':np.array([1,2,2,3,3,5])})
8 df.pivot_table(df, index=['A', 'B'], columns=['C'],
9 aggfunc=np.sum)
```

```
1 <tr>
2   <th></th>
3   <th></th>
4   <th colspan="3" halign="left">D</th>
5 </tr>
6 <tr>
7   <th></th>
8   <th>C</th>
9   <th>large</th>
10  <th>medium</th>
11  <th>small</th>
12 </tr>
13 <tr>
14  <th>A</th>
15  <th>B</th>
16  <th></th>
17  <th></th>
18  <th></th>
19 </tr>
```

```
1 <tr>
2   <th>bar</th>
3   <th>three</th>
4   <td>NaN</td>
5   <td>NaN</td>
6   <td>8.0</td>
7 </tr>
8 <tr>
9   <th rowspan="2" valign="top">foo</th>
10  <th>one</th>
11  <td>NaN</td>
12  <td>2.0</td>
13  <td>1.0</td>
14 </tr>
15 <tr>
16  <th>two</th>
17  <td>5.0</td>
18  <td>NaN</td>
19  <td>NaN</td>
20 </tr>
```

```
1 df =
2     pd.DataFrame({'A':np.array(['foo','foo','foo','foo','bar','bar']),
3
4         'B':np.array(['one','one','two','two','three','three']),
5
6         'C':np.array(['small','medium','large','large','small','small']),
7
8         'D':np.array([1,2,2,3,3,5]))}
5 df.groupby('A').agg(np.mean)
```

```
1 <tr style="text-align: right;">
2     <th></th>
3     <th>D</th>
4 </tr>
5 <tr>
6     <th>A</th>
7     <th></th>
8 </tr>
```

```
1 <tr>
2     <th>bar</th>
3     <td>4</td>
4 </tr>
5 <tr>
6     <th>foo</th>
7     <td>2</td>
8 </tr>
```

```
1 df.apply(np.mean, axis=0) # 对数据框df的每一列求平均值 axis: 0  
对列名（横着的）进行处理 1对索引（竖着的）进行处理
```

```
1 df =  
pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))  
2 df
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>A</th>  
4   <th>B</th>  
5   <th>C</th>  
6   <th>D</th>  
7   <th>E</th>  
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>0.800902</td>
4   <td>0.933676</td>
5   <td>0.461338</td>
6   <td>0.353398</td>
7   <td>0.057885</td>
8 </tr>
9 <tr>
10  <th>1</th>
11  <td>0.988580</td>
12  <td>0.681318</td>
13  <td>0.533361</td>
14  <td>0.486016</td>
15  <td>0.220004</td>
16 </tr>
17 <tr>
18  <th>2</th>
19  <td>0.695034</td>
20  <td>0.643920</td>
21  <td>0.694040</td>
22  <td>0.280063</td>
23  <td>0.641867</td>
24 </tr>
25 <tr>
26  <th>3</th>
27  <td>0.925290</td>
28  <td>0.084906</td>
29  <td>0.120247</td>
30  <td>0.880991</td>
31  <td>0.399596</td>
32 </tr>
33 <tr>
34  <th>4</th>
```

```
35 <td>0.697742</td>
36 <td>0.372860</td>
37 <td>0.881456</td>
38 <td>0.565627</td>
39 <td>0.272549</td>
40 </tr>
41 <tr>
42 <th>5</th>
43 <td>0.614245</td>
44 <td>0.658123</td>
45 <td>0.797487</td>
46 <td>0.609511</td>
47 <td>0.544633</td>
48 </tr>
49 <tr>
50 <th>6</th>
51 <td>0.153517</td>
52 <td>0.354870</td>
53 <td>0.910838</td>
54 <td>0.416895</td>
55 <td>0.098821</td>
56 </tr>
57 <tr>
58 <th>7</th>
59 <td>0.088223</td>
60 <td>0.501401</td>
61 <td>0.702754</td>
62 <td>0.334938</td>
63 <td>0.182708</td>
64 </tr>
65 <tr>
66 <th>8</th>
67 <td>0.737348</td>
68 <td>0.569340</td>
```

```
69 <td>0.291342</td>
70 <td>0.847058</td>
71 <td>0.193331</td>
72 </tr>
73 <tr>
74 <th>9</th>
75 <td>0.083915</td>
76 <td>0.396210</td>
77 <td>0.589415</td>
78 <td>0.806525</td>
79 <td>0.598841</td>
80 </tr>
```

```
1 df.apply(np.mean, axis=0)
```

```
1 A    0.578480
2 B    0.519662
3 C    0.598228
4 D    0.558102
5 E    0.321024
6 dtype: float64
```

数据的连接(join)与组合(combine)

```
1 df1.append(df2) # 在数据框df2的末尾添加数据框df1，其中df1和df2
                  的列数应该相等  列合并
```

```
1 df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],
2                     'B': ['B0', 'B1', 'B2', 'B3'],
3                     'C': ['C0', 'C1', 'C2', 'C3'],
4                     'D': ['D0', 'D1', 'D2', 'D3']},
5                     index=[0, 1, 2, 3])
6 df2 = pd.DataFrame({'A': ['A4', 'A5', 'A6', 'A7'],
7                     'B': ['B4', 'B5', 'B6', 'B7'],
8                     'C': ['C4', 'C5', 'C6', 'C7'],
9                     'D': ['D4', 'D5', 'D6', 'D7']},
10                    index=[4, 5, 6, 7])
11 df1.append(df2)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>A0</td>
4   <td>B0</td>
5   <td>C0</td>
6   <td>D0</td>
7 </tr>
8 <tr>
9   <th>1</th>
10  <td>A1</td>
11  <td>B1</td>
12  <td>C1</td>
13  <td>D1</td>
14 </tr>
15 <tr>
16   <th>2</th>
17   <td>A2</td>
18   <td>B2</td>
19   <td>C2</td>
20   <td>D2</td>
21 </tr>
22 <tr>
23   <th>3</th>
24   <td>A3</td>
25   <td>B3</td>
26   <td>C3</td>
27   <td>D3</td>
28 </tr>
29 <tr>
30   <th>4</th>
31   <td>A4</td>
32   <td>B4</td>
33   <td>C4</td>
34   <td>D4</td>
```

```
35 </tr>
36 <tr>
37   <th>5</th>
38   <td>A5</td>
39   <td>B5</td>
40   <td>C5</td>
41   <td>D5</td>
42 </tr>
43 <tr>
44   <th>6</th>
45   <td>A6</td>
46   <td>B6</td>
47   <td>C6</td>
48   <td>D6</td>
49 </tr>
50 <tr>
51   <th>7</th>
52   <td>A7</td>
53   <td>B7</td>
54   <td>C7</td>
55   <td>D7</td>
56 </tr>
```

```
1 pd.concat([df1, df2], axis=1) # 在数据框df1的列最后添加数据框
  df2,其中df1和df2的行数应该相等 # 中括号可以换成圆括号 # axis:
  0进行行合并 1进行列合并
```

```
1 df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],
2                     'B': ['B0', 'B1', 'B2', 'B3'],
3                     'C': ['C0', 'C1', 'C2', 'C3'],
4                     'D': ['D0', 'D1', 'D2', 'D3']},
5                     index=[0, 1, 2, 3])
6 df2 = pd.DataFrame({'A': ['A4', 'A5', 'A6', 'A7'],
7                     'B': ['B4', 'B5', 'B6', 'B7'],
8                     'C': ['C4', 'C5', 'C6', 'C7'],
9                     'D': ['D4', 'D5', 'D6', 'D7']},
10                    index=[4, 5, 6, 7])
11 pd.concat([df1, df2], axis=0)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>A0</td>
4   <td>B0</td>
5   <td>C0</td>
6   <td>D0</td>
7 </tr>
8 <tr>
9   <th>1</th>
10  <td>A1</td>
11  <td>B1</td>
12  <td>C1</td>
13  <td>D1</td>
14 </tr>
15 <tr>
16   <th>2</th>
17   <td>A2</td>
18   <td>B2</td>
19   <td>C2</td>
20   <td>D2</td>
21 </tr>
22 <tr>
23   <th>3</th>
24   <td>A3</td>
25   <td>B3</td>
26   <td>C3</td>
27   <td>D3</td>
28 </tr>
29 <tr>
30   <th>4</th>
31   <td>A4</td>
32   <td>B4</td>
33   <td>C4</td>
34   <td>D4</td>
```

```
35 </tr>
36 <tr>
37   <th>5</th>
38   <td>A5</td>
39   <td>B5</td>
40   <td>C5</td>
41   <td>D5</td>
42 </tr>
43 <tr>
44   <th>6</th>
45   <td>A6</td>
46   <td>B6</td>
47   <td>C6</td>
48   <td>D6</td>
49 </tr>
50 <tr>
51   <th>7</th>
52   <td>A7</td>
53   <td>B7</td>
54   <td>C7</td>
55   <td>D7</td>
56 </tr>
```

```
1 df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],
2                     'B': ['B0', 'B1', 'B2', 'B3'],
3                     'C': ['C0', 'C1', 'C2', 'C3'],
4                     'D': ['D0', 'D1', 'D2', 'D3']},
5                     index=[0, 1, 2, 3])
6 df2 = pd.DataFrame({'E': ['A4', 'A5', 'A6', 'A7'],
7                     'F': ['B4', 'B5', 'B6', 'B7'],
8                     'G': ['C4', 'C5', 'C6', 'C7'],
9                     'H': ['D4', 'D5', 'D6', 'D7']},
10                    index=[0, 1, 2, 3])
11 pd.concat((df1, df2), axis=1)
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>C</th>
6   <th>D</th>
7   <th>E</th>
8   <th>F</th>
9   <th>G</th>
10  <th>H</th>
11 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>A0</td>
4   <td>B0</td>
5   <td>C0</td>
6   <td>D0</td>
7   <td>A4</td>
8   <td>B4</td>
9   <td>C4</td>
10  <td>D4</td>
11 </tr>
12 <tr>
13  <th>1</th>
14  <td>A1</td>
15  <td>B1</td>
16  <td>C1</td>
17  <td>D1</td>
18  <td>A5</td>
19  <td>B5</td>
20  <td>C5</td>
21  <td>D5</td>
22 </tr>
23 <tr>
24  <th>2</th>
25  <td>A2</td>
26  <td>B2</td>
27  <td>C2</td>
28  <td>D2</td>
29  <td>A6</td>
30  <td>B6</td>
31  <td>C6</td>
32  <td>D6</td>
33 </tr>
34 <tr>
```

```
35 <th>3</th>
36 <td>A3</td>
37 <td>B3</td>
38 <td>C3</td>
39 <td>D3</td>
40 <td>A7</td>
41 <td>B7</td>
42 <td>C7</td>
43 <td>D7</td>
44 </tr>
```

```
1 df1.join(df2, on=col1, how='inner') # 对数据框df1和df2做内连接,
    其中连接的列为col1
```

```
1 df1 = pd.DataFrame({'A': ['A0', 'A1', 'A2', 'A3'],
2                     'B': ['B0', 'B1', 'B2', 'B3'],
3                     'key': ['K0', 'K1', 'K0', 'K1']})
4
5
6 df2 = pd.DataFrame({'C': ['C0', 'C1'],
7                     'D': ['D0', 'D1']},
8                     index=['K0', 'K1'])
9 df1.join(df2, on='key', how='inner')
```

```
1 <tr style="text-align: right;">
2   <th></th>
3   <th>A</th>
4   <th>B</th>
5   <th>key</th>
6   <th>C</th>
7   <th>D</th>
8 </tr>
```

```
1 <tr>
2   <th>0</th>
3   <td>A0</td>
4   <td>B0</td>
5   <td>K0</td>
6   <td>C0</td>
7   <td>D0</td>
8 </tr>
9 <tr>
10  <th>2</th>
11  <td>A2</td>
12  <td>B2</td>
13  <td>K0</td>
14  <td>C0</td>
15  <td>D0</td>
16 </tr>
17 <tr>
18  <th>1</th>
19  <td>A1</td>
20  <td>B1</td>
21  <td>K1</td>
22  <td>C1</td>
23  <td>D1</td>
24 </tr>
25 <tr>
26  <th>3</th>
27  <td>A3</td>
28  <td>B3</td>
29  <td>K1</td>
30  <td>C1</td>
31  <td>D1</td>
32 </tr>
```

数据的统计

```
1 df.mean() # 得到数据框df中每一列的平均值
```

```
1 df =  
pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))  
2 df.mean()
```

```
1 A    0.362158  
2 B    0.432248  
3 C    0.554478  
4 D    0.331155  
5 E    0.438283  
6 dtype: float64
```

```
1 df.corr() # 得到数据框df中每一列与其他列的相关系数
```

```
1 df.corr()
```

```
1 <tr style="text-align: right;">  
2   <th></th>  
3   <th>A</th>  
4   <th>B</th>  
5   <th>C</th>  
6   <th>D</th>  
7   <th>E</th>  
8 </tr>
```

```
1 <tr>
2   <th>A</th>
3   <td>1.000000</td>
4   <td>-0.167715</td>
5   <td>0.198216</td>
6   <td>0.036939</td>
7   <td>0.113714</td>
8 </tr>
9 <tr>
10  <th>B</th>
11  <td>-0.167715</td>
12  <td>1.000000</td>
13  <td>0.449789</td>
14  <td>0.015883</td>
15  <td>-0.236658</td>
16 </tr>
17 <tr>
18  <th>C</th>
19  <td>0.198216</td>
20  <td>0.449789</td>
21  <td>1.000000</td>
22  <td>-0.296943</td>
23  <td>0.386206</td>
24 </tr>
25 <tr>
26  <th>D</th>
27  <td>0.036939</td>
28  <td>0.015883</td>
29  <td>-0.296943</td>
30  <td>1.000000</td>
31  <td>-0.777327</td>
32 </tr>
33 <tr>
34  <th>E</th>
```

```
35 <td>0.113714</td>
36 <td>-0.236658</td>
37 <td>0.386206</td>
38 <td>-0.777327</td>
39 <td>1.000000</td>
40 </tr>
```

```
1 df.count() # 得到数据框df中每一列的非空值个数
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 # df.loc[0][0] = np.nan
4 df.iloc[0, 0] = np.nan
5 df.count()
```

```
1 A      9
2 B      10
3 C      10
4 D      10
5 E      10
6 dtype: int64
```

```
1 df.max() # 得到数据框df中每一列的最大值
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 df.max()
```

```
1 A    0.812708
2 B    0.886171
3 C    0.987035
4 D    0.977146
5 E    0.959625
6 dtype: float64
```

```
1 df.min() # 得到数据框df中每一列的最小值
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 df.min()
```

```
1 A    0.128560
2 B    0.135905
3 C    0.167476
4 D    0.137062
5 E    0.050306
6 dtype: float64
```

```
1 df.median() # 得到数据框df中每一列的中位数
```

```
1 df =
2 pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))
3 df.median()
```

```
1 A    0.409373
2 B    0.597418
3 C    0.678203
4 D    0.705762
5 E    0.519713
6 dtype: float64
```

```
1 | df.std() # 得到数据框df中每一列的标准差
```

```
1 | df =  
| pd.DataFrame(np.random.rand(10,5),columns=list('ABCDE'))  
2 | df.std()
```

```
1 | A      0.272847  
2 | B      0.254870  
3 | C      0.258956  
4 | D      0.301168  
5 | E      0.295753  
6 | dtype: float64
```