

05分组聚合

Pandas使用groupby实现分组统计

类似SQL:

```
select city,max(temperature) from city_weather group by city;
```

groupby: 先对数据分组，然后在每个分组上应用聚合函数、转换函数

本次演示:

- 1、分组使用聚合函数做数据统计
- 2、遍历groupby的结果理解执行流程
- 3、实例分组探索天气数据

```
1 # 定义一个df
2 df = pd.DataFrame({'A': ['foo', 'bar', 'foo', 'bar', 'foo',
3                           'bar', 'foo', 'foo'],
4                     'B': ['one', 'one', 'two', 'three', 'two',
5                           'two', 'one', 'three'],
6                     'C': np.random.randn(8),
7                     'D': np.random.randn(8)})
```

一、分组使用聚合函数做数据统计

1、单个列groupby，查询所有数据列的统计

```
1 df.groupby('A').sum()
2 # 可以看到输出结果中的'A'变成了索引，并且'B'列不是数值，所以被忽略
```

2、多个列groupby，查询所有数据列的统计

```
1 df.groupby(['A', 'B']).mean()
2 # ('A', 'B')成对变成了二级索引
3 df.groupby(['A', 'B'], as_index=False).mean()
4 # 去除二级索引
```

3、同时查看多种数据统计

```
1 df.groupby('A').agg([np.sum, np.mean, np.std])
2 # 列变成了多级索引
```

4、查看单列的结果数据统计

```
1 # 方法1: 预过滤, 性能更好
2 df.groupby('A')['C'].agg([np.sum, np.mean, np.std])
3 # 方法2
4 df.groupby('A').agg([np.sum, np.mean, np.std])['C']
```

5、不同列使用不同的聚合函数

```
1 df.groupby('A').agg({"C":np.sum, "D":np.mean})
```

二、遍历groupby的结果理解执行流程

for循环可以直接遍历每个group

1、遍历单个列聚合的分组

```
1 g = df.groupby('A')
2 g
3 # 获取单个分组的数据
4 g.get_group('bar')
```

2、遍历多个列聚合的分组

```
1 g = df.groupby(['A', 'B'])
2 for name, group in g:
3     print(name)
4     print(group)
5     print()
6 # 可以看到, 是一个2个元素的tuple, 代表不同的列
7 g.get_group(('foo', 'one'))
8 # 可以直接查询group后的某几列, 生成Series或者子DataFrame
9 g['C']
10 for name, group in g['C']:
11     print(name)
12     print(group)
13     print(type(group))
14     print()
15 # 其实所有的聚合统计, 都是在dataframe和series上进行的;
```

三、实例分组探索天气数据

```
1 fpath = "res/beijing_tianqi_2018.csv"
2 df = pd.read_csv(fpath)
3 # 替换掉温度的后缀℃
4 df.loc[:, "bwendu"] = df["bwendu"].str.replace("℃",
5     "").astype('int32')
6 df.loc[:, "ywendu"] = df["ywendu"].str.replace("℃",
7     "").astype('int32')
```

```

6 df.head()
7
8 # 1.新增一列为月份
9 df['month'] = df['ymd'].str[:7]
10 df.head()
11
12 # 2.查看每个月的最高温度
13 data = df.groupby('month')['bwendu'].max()
14 data
15 type(data)
16 data.plot() # 画图
17
18 # 3.查看每个月的最高温度、最低温度、平均空气质量指数
19 group_data = df.groupby('month').agg({"bwendu":np.max,
20 "ywendu":np.min, "aqi":np.mean})
21 group_data
22 group_data.plot()

```

四、分组聚合后不同列的数据统计

电影评分数据集 (UserID, MovieID, Rating, Timestamp)

聚合后单列-单指标统计：每个MovieID的平均评分

```
df.groupby("MovieID")["Rating"].mean()
```

聚合后单列-多指标统计：每个MovieID的最高评分、最低评分、平均评分

```
df.groupby("MovieID")["Rating"].agg(mean="mean", max="max", min=np.min)
df.groupby("MovieID").agg({"Rating":['mean', 'max', np.min]})
```

聚合后多列-多指标统计：每个MovieID的评分人数，最高评分、最低评分、平均评分

```
df.groupby("MovieID").agg( rating_mean=("Rating", "mean"), user_count=
("UserID", lambda x : x.nunique()))
df.groupby("MovieID").agg( {"Rating": ['mean', 'min', 'max'], "UserID":
lambda x :x.nunique()})
df.groupby("MovieID").apply( lambda x: pd.Series( {"min":
x["Rating"].min(), "mean": x["Rating"].mean()}))
```

记忆：agg(新列名=函数)、agg(新列名=(原列名, 函数))、agg({"原列名": 函数/列表})

agg函数的两种形式，等号代表“把结果赋值给新列”，字典/元组代表“对这个列运用这些函数”

官网文档：<https://pandas.pydata.org/pandas-docs/version/0.23.4/generated/pandas.core.groupby.DataFrameGroupBy.agg.html>

1.读取数据

```

1 df = pd.read_csv(
2     "res/ratings.dat",
3     sep="::",
4     engine='python',
5     names="UserID::MovieID::Rating::Timestamp".split("::")
6 )

```

2.聚合后单列-单指标统计

```

1 # 每个MovieID的平均评分
2 result = df.groupby("MovieID")["Rating"].mean()
3 result.head()
4 type(result)

```

3.聚合后单列-多指标统计

每个MovieID的最高评分、最低评分、平均评分

```

1 # 方法1: agg函数传入多个结果列名=函数名形式
2 result = df.groupby("MovieID")["Rating"].agg(
3     mean="mean", max="max", min=np.min
4 )
5 result.head()
6
7 # 方法2: agg函数传入字典, key是column名, value是函数列表
8 # 每个MovieID的最高评分、最低评分、平均评分
9 result = df.groupby("MovieID").agg(
10     {"Rating": ['mean', 'max', np.min]}
11 )
12 result.head()
13 # 去除二级索引
14 result.columns = ['age_mean', 'age_min', 'age_max']
15 result.head()

```

4.聚合后多列-多指标统计

每个MovieID的评分人数, 最高评分、最低评分、平均评分

```

1 # 方法1: agg函数传入字典, key是原列名, value是原列名和函数元组
2 result = df.groupby("MovieID").agg(
3     rating_mean=("Rating", "mean"),
4     rating_min=("Rating", "min"),
5     rating_max=("Rating", "max"),
6     user_count=("UserID", lambda x : x.nunique())
7 )
8 result.head()
9
10 # 方法2: agg函数传入字典, key是原列名, value是函数列表
11 result = df.groupby("MovieID").agg(

```

```
12     {
13         "Rating": ['mean', 'min', 'max'],
14         "UserID": lambda x :x.nunique()
15     }
16 )
17 result.head() # 结果是二级索引
18 # 去除二级索引
19 result.columns = ["rating_mean",
20                  "rating_min", "rating_max", "user_count"]
21 result.head()
22 # 方法3: 使用groupby之后apply对每个子df单独统计
23 def agg_func(x):
24     """注意, 这个x是子DF"""
25
26     # 这个Series会变成一行, 字典KEY是列名
27     return pd.Series({
28         "rating_mean": x["Rating"].mean(),
29         "rating_min": x["Rating"].min(),
30         "rating_max": x["Rating"].max(),
31         "user_count": x["UserID"].nunique()
32     })
33
34 result = df.groupby("MovieID").apply(agg_func)
35 result.head()
36
```