6.13 数据的累加与统计操作

问题¶

你需要处理一个很大的数据集并需要计算数据总和或其他统计量。

解决方案¶

对于任何涉及到统计、时间序列以及其他相关技术的数据分析问题,都可以考虑使用 Pandas库。

为了让你先体验下,下面是一个使用Pandas来分析芝加哥城市的 <u>老鼠和啮齿类动物数据库</u> 的例子。 在我写这篇文章的时候,这个数据库是一个拥有大概74,000行数据的CSV文件。

```
>>> import pandas
>>> # Read a CSV file, skipping last line
>>> rats = pandas.read csv('rats.csv', skip footer=1)
>>> rats
<class 'pandas.core.frame.DataFrame'>
Int64Index: 74055 entries, 0 to 74054
Data columns:
Creation Date 74055 non-null values
Status 74055 non-null values
Completion Date 72154 non-null values
Service Request Number 74055 non-null values
Type of Service Request 74055 non-null values
Number of Premises Baited 65804 non-null values
Number of Premises with Garbage 65600 non-null values
Number of Premises with Rats 65752 non-null values
Current Activity 66041 non-null values
Most Recent Action 66023 non-null values
Street Address 74055 non-null values
ZIP Code 73584 non-null values
X Coordinate 74043 non-null values
Y Coordinate 74043 non-null values
Ward 74044 non-null values
Police District 74044 non-null values
Community Area 74044 non-null values
Latitude 74043 non-null values
Longitude 74043 non-null values
Location 74043 non-null values
dtypes: float64(11), object(9)
>>> # Investigate range of values for a certain field
>>> rats['Current Activity'].unique()
array([nan, Dispatch Crew, Request Sanitation Inspector], dtype=object)
>>> # Filter the data
>>> crew_dispatched = rats[rats['Current Activity'] == 'Dispatch Crew']
>>> len(crew_dispatched)
65676
>>> # Find 10 most rat-infested ZIP codes in Chicago
>>> crew dispatched['ZIP Code'].value counts()[:10]
60647 3837
60618 3530
60614 3284
60629 3251
60636 2801
60657 2465
60641 2238
60609 2206
60651 2152
60632 2071
>>>
```

```
>>> # Group by completion date
>>> dates = crew dispatched.groupby('Completion Date')
<pandas.core.groupby.DataFrameGroupBy object at 0x10d0a2a10>
>>> len(dates)
472
>>>
>>> # Determine counts on each day
>>> date counts = dates.size()
>>> date counts[0:10]
Completion Date
01/03/2011 4
01/03/2012 125
01/04/2011 54
01/04/2012 38
01/05/2011 78
01/05/2012 100
01/06/2011 100
01/06/2012 58
01/07/2011 1
01/09/2012 12
>>>
>>> # Sort the counts
>>> date counts.sort()
>>> date counts[-10:]
Completion Date
10/12/2012 313
10/21/2011 314
09/20/2011 316
10/26/2011 319
02/22/2011 325
10/26/2012 333
03/17/2011 336
10/13/2011 378
10/14/2011 391
10/07/2011 457
>>>
```

嗯,看样子2011年10月7日对老鼠们来说是个很忙碌的日子啊! ^ ^

讨论¶

Pandas是一个拥有很多特性的大型函数库,我在这里不可能介绍完。 但是只要你需要去分析大型数据集合、对数据分组、计算各种统计量或其他类似任务的话,这个函数库真的值得你去看一看。