

上机作业 1：正则表达式应用

088 于海鑫

2017211240

name1e5s@bupt.edu.cn

版本：8

更新：March 23, 2020

1 实验要求

从因特网上搜索 Web 页，用 `wget` 获取网页，处理网页 `html` 文本数据，从中提取出当前时间点北京各监测站的 PM2.5 浓度，输出如下 CSV 格式数据。

```
1 2020-03-09 13:00:00, 海淀区万柳, 73
2 2020-03-09 13:00:00, 昌平镇, 67
3 2020-03-09 13:00:00, 奥体中心, 66
4 2020-03-09 13:00:00, 海淀区万柳, 73
5 2020-03-09 13:00:00, 昌平镇, 73
6 2020-03-09 13:00:00, 奥体中心, 75
```

撰写实验报告。

2 实验过程

2.1 获取数据

在 `bing` 上搜索“北京各监测站的 PM2.5 浓度”并筛选，最终我们选定了最为简单的 `http://www.86pm25.com/city/beijing.html` 作为筛选的目标。在这个网页中，我们要进行处理的 HTML 片段如下：

2.2.1 提取表格

我们首先要做的就是把这一部分表格从整个 `HTML` 文本中筛选出来，经过观察可以发现该网页的表格十分有规律，有且只有表格部分的开头为 `<tr><td>`，以此为切入点，我们使用 `grep` 配合正则表达式抓取表格片段。

```
1 bash$ grep "^<tr><td>"
```

2.3 删除 TAG

`HTML` 作为富文本文档，在转换为我们需要的格式前必须要做的就是删除多余的 `TAG` 部分给清除掉，这里我们使用 `sed` 来解决问题，代码如下：

```
1 bash$ sed -E 's/(<[^>]+>)+/~ /g' | sed 's/~ / /g'
```

这一段代码的作用时先把全部的 `TAG` 转换为单个的波浪号字符，之后再把该字符转换为空格。此时的输出如下：

```
1 name1e5s@sumeru:~$ wget -q -O - http://www.86pm25.com/city/beijing.html | grep "^<tr><td>" | sed -E 's/(<[^>]+>)+/~ /g' | sed 's/~ / /g'
2 '
3 奥体中心 28 13µg/m³ 21µg/m³
4 昌平镇 28 7µg/m³ 10µg/m³
5 定陵 29 7µg/m³ 12µg/m³
6 东四 30 17µg/m³ 28µg/m³
7 古城 29 13µg/m³ 17µg/m³
8 官园 33 11µg/m³ 18µg/m³
9 海淀区万柳 30 10µg/m³ 10µg/m³
10 农展馆 34 17µg/m³ 34µg/m³
11 顺义新城 32 22µg/m³ 28µg/m³
12 天坛 37 17µg/m³ 37µg/m³
13 万寿西宫 24 13µg/m³ 24µg/m³
```

2.4 删除单位

`µg/m³` 作为不需要的单位，使用 `sed` 剔除掉。

```
1 bash$ sed 's/%µg/m³%%g'
```

此时的输出如下：

```
1 name1e5s@sumeru:~$ wget -q -O - http://www.86pm25.com/city/beijing.  
  html | grep "^<tr><td>" | sed -E 's/(<[^>]+>)+/~ /g' | sed 's/~ / /g  
  ' | sed 's/%μg/m³%%g'  
2 奥体中心 28    13 21  
3 昌平镇 28     7 10  
4 定陵 29     7 12  
5 东四 30    17 28  
6 古城 29    13 17  
7 官园 33    11 18  
8 海淀区万柳 30    10 10  
9 农展馆 34    17 34  
10 顺义新城 32    22 28  
11 天坛 37    17 37  
12 万寿西宫 24    13 24
```

2.5 打印结果

在获取到必要的信息之后，我们可以使用 `awk` 输出结果。

```
1 bash$ awk '{print datenow "," $1 "," $3}' datenow="`date "+%Y-%m-%d %  
  H:%M:%S"``"
```

此时的输出如下：

```
1 name1e5s@sumeru:~$ wget -q -O - http://www.86pm25.com/city/beijing.  
  html | grep "^<tr><td>" | sed -E 's/(<[^>]+>)+/~ /g' | sed 's/~ / /g  
  ' | sed 's/%μg/m³%%g'| awk '{print datenow "," $1 "," $3}' datenow  
  ="`date "+%Y-%m-%d %H:%M:%S"``"  
2 2020-03-23 14:35:43,奥体中心,13  
3 2020-03-23 14:35:43,昌平镇,7  
4 2020-03-23 14:35:43,定陵,7  
5 2020-03-23 14:35:43,东四,17  
6 2020-03-23 14:35:43,古城,13  
7 2020-03-23 14:35:43,官园,11  
8 2020-03-23 14:35:43,海淀区万柳,10  
9 2020-03-23 14:35:43,农展馆,17  
10 2020-03-23 14:35:43,顺义新城,22
```

- | | |
|----|-------------------------------|
| 11 | 2020-03-23 14:35:43, 天坛, 17 |
| 12 | 2020-03-23 14:35:43, 万寿西宫, 13 |

符合要求。