

# Linux作业一

- 学号：2018211582
- 班级：2018211314
- 姓名：李志毅

## 一、实验准备

搜索含有北京PM2.5的网站，找到网址为<http://www.86pm25.com/city/beijing.html>，网址含有的元素如下：

各监测站点实时数据				
监测站点	AQI	污染等级	PM2.5浓度	PM10浓度
奥体中心	117	<div></div> 轻度污染	20µg/m³	183µg/m³
昌平镇	135	<div></div> 轻度污染	14µg/m³	219µg/m³
大兴旧宫	126	<div></div> 轻度污染	20µg/m³	201µg/m³
走马(对照点)	125	<div></div> 轻度污染	18µg/m³	200µg/m³
东四	141	<div></div> 轻度污染	16µg/m³	232µg/m³
房山燕山	97	<div></div> 良	16µg/m³	144µg/m³
丰台小屯	135	<div></div> 轻度污染	25µg/m³	219µg/m³
丰台云岗	122	<div></div> 轻度污染	17µg/m³	194µg/m³
古城	125	<div></div> 轻度污染	19µg/m³	199µg/m³
官园	113	<div></div> 轻度污染	8µg/m³	176µg/m³
海淀万柳	124	<div></div> 轻度污染	21µg/m³	198µg/m³
怀柔新城	161	<div></div> 中度污染	17µg/m³	271µg/m³
门头沟三家店	144	<div></div> 轻度污染	21µg/m³	238µg/m³
密云新城	170	<div></div> 中度污染	18µg/m³	289µg/m³
密云镇	158	<div></div> 中度污染	18µg/m³	265µg/m³
农展馆	124	<div></div> 轻度污染	12µg/m³	197µg/m³
平谷新城	166	<div></div> 中度污染	20µg/m³	281µg/m³
顺义新城	167	<div></div> 中度污染	20µg/m³	283µg/m³
天坛	116	<div></div> 轻度污染	20µg/m³	182µg/m³
通州东关	163	<div></div> 中度污染	21µg/m³	275µg/m³
万寿西宫	112	<div></div> 轻度污染	19µg/m³	173µg/m³
延庆石河营	138	<div></div> 轻度污染	24µg/m³	225µg/m³
延庆夏都	147	<div></div> 轻度污染	18µg/m³	244µg/m³

使用F12观察网址HTML5结构，发现我们需要的数据存放在tr和td标签中，含有监测地点，PM2.5指数，在网站上部分标签中含有需要的时间数据



## 1. 替换标签

<tr><td>奥体中心</td><td>238</td><td><img src='../images/wurandengji/zhongdu.gif' /></td><td>188μg/m³</td><td>298μg/m³</td></tr>  
 <tr><td>昌平镇</td><td>221</td><td><img src='../images/wurandengji/zhongdu.gif' /></td><td>171μg/m³</td><td>231μg/m³</td></tr>  
 <tr><td>大兴旧宫</td><td>237</td><td><img src='../images/wurandengji/zhongdu.gif' /></td><td>187μg/m³</td><td>231μg/m³</td></tr>  
 <tr><td>定陵(对照点)</td><td>219</td><td><img src='../images/wurandengji/zhongdu.gif' /></td><td>169μg/m³</td><td>210μg/m³</td></tr>  
 <tr><td>东四</td><td>230</td><td><img src='../images/wurandengji/zhongdu.gif' /></td><td>180μg/m³</td><td>224μg/m³</td></tr>

```
sed -e 's/<[^\<>]*>/ /g'
```

```
a582@Ubuntu:~$ cat beijing.html | sed -e 's/<[^<>]*>/ /g'
```

各监测站点实时数据				
监测站点	AQI	污染等级	PM2.5浓度	PM10浓度
奥体中心 238	188 $\mu\text{g}/\text{m}^3$	298 $\mu\text{g}/\text{m}^3$		
昌平镇 221	171 $\mu\text{g}/\text{m}^3$	231 $\mu\text{g}/\text{m}^3$		
大兴旧宫 237	187 $\mu\text{g}/\text{m}^3$	231 $\mu\text{g}/\text{m}^3$		
定陵(对照点) 219	169 $\mu\text{g}/\text{m}^3$	210 $\mu\text{g}/\text{m}^3$		
东四 230	180 $\mu\text{g}/\text{m}^3$	224 $\mu\text{g}/\text{m}^3$		
房山燕山 172	130 $\mu\text{g}/\text{m}^3$	173 $\mu\text{g}/\text{m}^3$		
丰台小屯 209	159 $\mu\text{g}/\text{m}^3$	249 $\mu\text{g}/\text{m}^3$		
丰台云岗 170	129 $\mu\text{g}/\text{m}^3$	230 $\mu\text{g}/\text{m}^3$		
古城 199	149 $\mu\text{g}/\text{m}^3$	227 $\mu\text{g}/\text{m}^3$		
官园 216	166 $\mu\text{g}/\text{m}^3$	237 $\mu\text{g}/\text{m}^3$		
海淀万柳 217	167 $\mu\text{g}/\text{m}^3$	248 $\mu\text{g}/\text{m}^3$		
怀柔新城 195	146 $\mu\text{g}/\text{m}^3$	177 $\mu\text{g}/\text{m}^3$		
怀柔镇 203	153 $\mu\text{g}/\text{m}^3$	172 $\mu\text{g}/\text{m}^3$		
门头沟三家店 188	141 $\mu\text{g}/\text{m}^3$	212 $\mu\text{g}/\text{m}^3$		
密云新城 211	161 $\mu\text{g}/\text{m}^3$	204 $\mu\text{g}/\text{m}^3$		
密云镇 204	154 $\mu\text{g}/\text{m}^3$	196 $\mu\text{g}/\text{m}^3$		
农展馆 235	185 $\mu\text{g}/\text{m}^3$	276 $\mu\text{g}/\text{m}^3$		
平谷新城 130	99 $\mu\text{g}/\text{m}^3$	141 $\mu\text{g}/\text{m}^3$		
顺义新城 210	160 $\mu\text{g}/\text{m}^3$	226 $\mu\text{g}/\text{m}^3$		
通州东关 223	173 $\mu\text{g}/\text{m}^3$	225 $\mu\text{g}/\text{m}^3$		
延庆石河营 150	115 $\mu\text{g}/\text{m}^3$	143 $\mu\text{g}/\text{m}^3$		
延庆夏都 176	133 $\mu\text{g}/\text{m}^3$	149 $\mu\text{g}/\text{m}^3$		

北京实时空气质量指数  
更新：2021年03月26日 11时

```
cat beijing.html | sed -e 's/<[^<>]*>/ /g' -e 's/[年月]/-/g'
```

北京实时空气质量指数  
更新: 2021-03-26日 11时

### 3.处理PM2.5数据并打印

之后我们要将含有时间的这一行和包含PM2.5数据的行筛选出来, 我们可以编写 `awk` 文件 `flow.awk`, 将含有文本 "更新: " 的一行中的第一个数据和第二个 数据分别赋值给变量`date`和`time`, 将含有  $\text{m}^3$  的一行中取第一个数据 监测地点 和第三个数据 `PM2.5`指数, 使用 `printf` 语句打印数据, `flow.awk` 文件内容如下:

```
/更新:/{date=$1;time=$2;}  
/m³/{  
    printf("%s %s,%s,%s\n",date,time,$1,$3);  
}
```

```
a582@Ubuntu:~$ cat beijing.html | sed -e 's/<[^<>]*>/ /g' -e 's/[年月]/-/g' | awk -f flow.awk | more  
更新: 2021-03-26日 11时,奥体中心,188μg/m³  
更新: 2021-03-26日 11时,昌平镇,171μg/m³  
更新: 2021-03-26日 11时,大兴旧宫,187μg/m³  
更新: 2021-03-26日 11时,定陵(对照点),169μg/m³  
更新: 2021-03-26日 11时,东四,180μg/m³  
更新: 2021-03-26日 11时,房山燕山,130μg/m³  
更新: 2021-03-26日 11时,丰台小屯,159μg/m³  
更新: 2021-03-26日 11时,丰台云岗,129μg/m³  
更新: 2021-03-26日 11时,古城,149μg/m³  
更新: 2021-03-26日 11时,官园,166μg/m³  
更新: 2021-03-26日 11时,海淀万柳,167μg/m³  
更新: 2021-03-26日 11时,怀柔新城,146μg/m³  
更新: 2021-03-26日 11时,怀柔镇,153μg/m³  
更新: 2021-03-26日 11时,门头沟三家店,141μg/m³  
更新: 2021-03-26日 11时,密云新城,161μg/m³  
更新: 2021-03-26日 11时,密云镇,154μg/m³  
更新: 2021-03-26日 11时,农展馆,185μg/m³  
更新: 2021-03-26日 11时,平谷新城,99μg/m³  
更新: 2021-03-26日 11时,顺义新城,160μg/m³  
更新: 2021-03-26日 11时,通州东关,173μg/m³  
更新: 2021-03-26日 11时,延庆石河营,115μg/m³  
更新: 2021-03-26日 11时,延庆夏都,133μg/m³
```

### 4.整理最终数据

最后我们整理一下这个数据, 将 "更新: " 替换成空, 将 "时" 替换成 ":00:00", 将  $\text{ug}/\text{m}^3$  替换成空, 将 "日" 替换成空, 使用`sed`语句:

```
sed -e 's/μg.m³/ /g' -e 's/[更新: 日]/ /g' -e 's/时/:00:00/g' | more
```

最终的完整的命令如下:

```
cat beijing.html | sed -e 's/<[^<>]*>/ /g' -e 's/[年月]/-/g' | awk -f flow.awk |  
sed -e 's/μg.m³?/ /g' -e 's/[更新: 日]/ /g' -e 's/时/:00:00/g' | more
```

最终执行后的截图如下, 包含监测时间、监测地点和监测地点的PM2.5指数, 与实验要求一致

```
a582@ubuntu:~$ cat beijing.html | sed -e 's/<[^>]*>/ /g' -e 's/[年月日-/:/g'| awk -f flow.awk | sed -e 's/μg.m³/ /g' -e 's/[更新: 日]/g' -e 's/时/:00:00/g' | more
2021-03-26 11:00:00,奥体中心,188
2021-03-26 11:00:00,昌平镇,171
2021-03-26 11:00:00,大兴旧宫,187
2021-03-26 11:00:00,定陵(对焦点),169
2021-03-26 11:00:00,东四,180
2021-03-26 11:00:00,房山燕山,130
2021-03-26 11:00:00,丰台小屯,159
2021-03-26 11:00:00,丰台云岗,129
2021-03-26 11:00:00,古城,149
2021-03-26 11:00:00,官园,166
2021-03-26 11:00:00,海淀万柳,167
2021-03-26 11:00:00,怀柔城,146
2021-03-26 11:00:00,怀柔镇,153
2021-03-26 11:00:00,门头沟三家店,141
2021-03-26 11:00:00,密云城,161
2021-03-26 11:00:00,密云镇,154
2021-03-26 11:00:00,农展馆,195
2021-03-26 11:00:00,平谷城,99
2021-03-26 11:00:00,顺义城,160
2021-03-26 11:00:00,通州东关,173
2021-03-26 11:00:00,延庆古河套,115
2021-03-26 11:00:00,延庆夏都,133
```

### 三、实验问题

在本次实验中我遇到了以下几个问题，并最终解决了这些问题

#### 1.网址选择问题

包含北京PM2.5数据的网址有很多，不同的网站的HTML5结构不同，数据包含在不同的标签内，因此选择合适的网址可以进行更简便的正则匹配，最后在综合考虑多方面因素后，我选择了本次实验的网址，他包含本次实验所需的所有数据，并且网站结构简单，没有复杂的标签，数据集中，适合于本次实验的正则匹配实验。

#### 2.sed指令使用不熟练

在刚开始实验的时候，sed指令只在课上学习过，并跟着MOOC上的案例实验过，并没有亲自动手在Linux系统上使用过，因此在最开始处理数据时，对于sed用于替换的用法不熟悉，容易忘记单引号等。当需要复杂的替换语句时，更需要多次观看MOOC上的内容(例如第一个sed语句，对HTML5中的标签替换)，来写出最简洁的正则匹配语句。

#### 3.awk文件使用不熟练

在进行必要的一些替换后，需要编写awk脚本用户选择并打印出我们需要的数据的行，在刚开始编写时，需要选择该行中标志性的特殊词，以此定位到某一行，例如需要时间这个数据，在H5中搜寻后发现，该行中包含“更新”两个特殊的语句，可以以此为切入点拿到数据，但拿到的该行第一列数据中除了时间外还包含有“更新”两个字，同时还有多余的“日”字，这些都在awk脚本执行后再执行sed语句处理，对于包含PM2.5数据的行，我们以m³为切入点，打印每一行，并在最后使用sed语句将m³剔除。编写awk文件时，找到所需行中的特殊字符或字符串很重要。

#### 4.命令使用不够细节

在第一次编写完成并测试时，我发现我筛选出的数据都显示在一行中，多个PM2.5数据并没有换行，一开始我认为可能是元HTML5文件的tr标签之间并没有换行，我使用sed语句为每个后加入了换行符，发现最终的数据依旧没有换行，经过仔细的梳理，发现了原来是awk脚本中printf语句最后没有加入换行符\n，这是由于实验不够细致导致的，也为我以后进行实验敲响了警钟，要细致的处理数据，耐心的找寻出现的问题并解决。

### 四、实验总结

本次实验通过搜寻北京PM2.5Web网页并处理该网页上的HTML5数据，练习了正则匹配表达式的应用，本次实验让我正则表达式的运用能力得到了进一步的提升，对于sed语句和awk脚本文件的使用变得更加的熟练，对于在Linux系统上使用正则表达式处理数据有了更深的体会，从最开始使用指令都不熟悉到最后可以熟练根据文本内容分析出应该使用的语句，这次实验使我收获颇丰。