

# Awk实战案例精讲

我：骏马金龙

博客：骏马金龙 [www.junmajinlong.com](http://www.junmajinlong.com)

Shell技术交流QQ群：921383787

## 插入几个新字段

在"a b c d"的b后面插入3个字段 e f g 。

```
1 echo a b c d|awk '{ $3="e f g " $3 }1'
```

## 格式化空白

移除每行的前缀、后缀空白，并将各部分左对齐。

```
1      aaaa      bbb      ccc
2      bbb      aaa ccc
3 ddd      fff      eee gg hh ii jj
```

```
1 awk 'BEGIN{OFS="\t"}{$1=$1;print}' a.txt
```

执行结果：

```
1 aaaa      bbb      ccc
2 bbb      aaa      ccc
3 ddd      fff      eee      gg      hh      ii      jj
```

## 筛选IPv4地址

从ifconfig命令的结果中筛选出除了lo网卡外的所有IPv4地址。

```
1  ## 1.法一:
2  ifconfig | awk '/inet / && !($2 ~ /^127/){print $2}'
3
4  # 按段落读取
5  ## 2.法二:
6  ifconfig | awk 'BEGIN{RS=""}!/lo/{print $6}'
7
8  ## 3.法三:
9  ifconfig |\
10 awk '
11     BEGIN{RS="";FS="\n"}
12     !/lo/{ $0=$2;FS=" ";$0=$0;print $2}
13 '
```

## 读取.ini配置文件中的某段

```
1 [base]
2 name=os_repo
```

```

3  baseurl=https://xxx/centos/$releasever/os/$basearch
4  gpgcheck=0
5
6  enable=1
7
8  [mysql]
9  name=mysql_repo
10 baseurl=https://xxx/mysql-repo/yum/mysql-5.7-community/el/$releasever/$basearch
11
12 gpgcheck=0
13 enable=1
14
15 [epel]
16 name=epel_repo
17 baseurl=https://xxx/epel/$releasever/$basearch
18 gpgcheck=0
19 enable=1
20 [percona]
21 name=percona_repo
22 baseurl = https://xxx/percona/release/$releasever/RPMS/$basearch
23 enabled = 1
24 gpgcheck = 0

```

```

1  awk '
2      BEGIN{RS="" } # 按段落
3      /\[mysql\]/{
4          print;
5          while( (getline)>0 ){
6              if(/\[.*\]/){
7                  exit
8              }
9              print
10         }
11     }' a.txt

```

## 根据某字段去重

去掉 `uid=xxx` 重复的行。

```

1  2019-01-13_12:00_index?uid=123
2  2019-01-13_13:00_index?uid=123
3  2019-01-13_14:00_index?uid=333
4  2019-01-13_15:00_index?uid=9710
5  2019-01-14_12:00_index?uid=123
6  2019-01-14_13:00_index?uid=123
7  2019-01-15_14:00_index?uid=333
8  2019-01-16_15:00_index?uid=9710

```

```
1  awk -F"?" ' !arr[$2]++{print}' a.txt
```

结果:

```
1 2019-01-13_12:00_index?uid=123
2 2019-01-13_14:00_index?uid=333
3 2019-01-13_15:00_index?uid=9710
```

## 次数统计

---

```
1 portmapper
2 portmapper
3 portmapper
4 portmapper
5 portmapper
6 portmapper
7 status
8 status
9 mountd
10 mountd
11 mountd
12 mountd
13 mountd
14 mountd
15 nfs
16 nfs
17 nfs_acl
18 nfs
19 nfs
20 nfs_acl
21 nlockmgr
22 nlockmgr
23 nlockmgr
24 nlockmgr
25 nlockmgr
```

```
1 awk '
2 {arr[$1]++}
3 END{
4     OFS="\t";
5     for(idx in arr){printf arr[idx],idx}
6 }
7 ' a.txt
```

## 统计TCP连接状态数量

---

```

1  $ netstat -tnap
2  Proto Recv-Q Send-Q Local Address   Foreign Address  State           PID/Program name
3  tcp        0      0 0.0.0.0:22       0.0.0.0:*        LISTEN          1139/sshd
4  tcp        0      0 127.0.0.1:25     0.0.0.0:*        LISTEN          2285/master
5  tcp        0     96 192.168.2.17:22  192.168.2.1:2468 ESTABLISHED    87463/sshd: root@pt
6  tcp        0      0 192.168.2017:22  192.168.201:5821 ESTABLISHED    89359/sshd: root@no
7  tcp6       0      0 :::3306          :::*             LISTEN          2289/mysqld
8  tcp6       0      0 :::22           :::*             LISTEN          1139/sshd
9  tcp6       0      0 :::1:25         :::*             LISTEN          2285/master

```

统计得到的结果：

```

1  5: LISTEN
2  2: ESTABLISHED

```

```

1  netstat -tnap | \
2  awk '
3      /^tcp/{
4          arr[$6]++
5      }
6      END{
7          for(state in arr){
8              print arr[state] ": " state
9          }
10     }
11 '

```

一行式：

```

1  netstat -tna | awk '/^tcp/{arr[$6]++}END{for(state in arr){print arr[state] ": " state}}'
2  netstat -tna | /usr/bin/grep 'tcp' | awk '{print $6}' | sort | uniq -c

```

## 统计日志中各IP访问非200状态码的次数

日志示例数据：

```

1  111.202.100.141 - - [2019-11-07T03:11:02+08:00] "GET /robots.txt HTTP/1.1" 301 169

```

统计非200状态码的IP，并取次数最多的前10个IP。

```

1  # 法一
2  awk '
3      $8!=200{arr[$1]++}
4      END{
5          for(i in arr){print arr[i],i}
6      }
7  ' access.log | sort -k1nr | head -n 10
8
9  # 法二：
10 awk '
11     $8!=200{arr[$1]++}

```

```

12     END{
13         PROCINFO["sorted_in"]="@val_num_desc";
14         for(i in arr){
15             if(cnt++==10){exit}
16             print arr[i],i
17         }
18     }' access.log

```

## 统计独立IP

url 访问IP 访问时间 访问人

```

1  a.com.cn|202.109.134.23|2015-11-20 20:34:43|guest
2  b.com.cn|202.109.134.23|2015-11-20 20:34:48|guest
3  c.com.cn|202.109.134.24|2015-11-20 20:34:48|guest
4  a.com.cn|202.109.134.23|2015-11-20 20:34:43|guest
5  a.com.cn|202.109.134.24|2015-11-20 20:34:43|guest
6  b.com.cn|202.109.134.25|2015-11-20 20:34:48|guest

```

需求：统计每个URL的独立访问IP有多少个(去重)，并且要为每个URL保存一个对应的文件，得到的结果类似：

```

1  a.com.cn  2
2  b.com.cn  2
3  c.com.cn  1

```

并且有三个对应的文件：

```

1  a.com.cn.txt
2  b.com.cn.txt
3  c.com.cn.txt

```

代码：

```

1  BEGIN{
2      FS="|"
3  }
4
5  !arr[$1,$2]++{
6      arr1[$1]++
7  }
8
9  END{
10     for(i in arr1){
11         print i,arr1[i] >("i.txt")
12     }
13 }

```

## 处理字段缺失的数据

	ID	name	gender	age	email	phone
1	1	Bob	male	28	abc@qq.com	18023394012
2	2	Alice	female	24	def@gmail.com	18084925203
3	3	Tony	male	21		17048792503
4	4	Kevin	male	21	bbb@189.com	17023929033
5	5	Alex	male	18	ccc@xyz.com	18185904230
6	6	Andy	female		ddd@139.com	18923902352
7	7	Jerry	female	25	exdsa@189.com	18785234906
8	8	Peter	male	20	bax@qq.com	17729348758
9	9	Steven		23	bc@sohu.com	15947893212
10	10	Bruce	female	27	bcbd@139.com	13942943905

当字段缺失时，直接使用FS划分字段来处理会非常棘手。gawk为了解决这种特殊需求，提供了FIELDWIDTHS变量。

FIELDWIDTHS可以按照字符数量划分字段。

```
1 awk '{print $4}' FIELDWIDTHS="2 2:6 2:6 2:3 2:13 2:11" a.txt
```

## 处理字段中包含了字段分隔符的数据

下面是CSV文件中的一行，该CSV文件以逗号分隔各个字段。

```
1 Robbins,Arnold,"1234 A Pretty Street, NE",MyTown,MyState,12345-6789,USA
```

需求：取得第三个字段"1234 A Pretty Street, NE"。

当字段中包含了字段分隔符时，直接使用FS划分字段来处理会非常棘手。gawk为了解决这种特殊需求，提供了FPAT变量。

FPAT可以收集正则匹配的结果，并将它们保存在各个字段中。（就像grep匹配成功的部分会加颜色显示，而使用FPAT划分字段，则是将匹配成功的部分保存在字段 `$1 $2 $3...` 中）。

```
1 echo 'Robbins,Arnold,"1234 A Pretty Street, NE",MyTown,MyState,12345-6789,USA' | \
2 awk 'BEGIN{FPAT="[^,]+|\".*\""}{print $1,$3}'
```

## 取字段中指定字符数量

```
1 16 001agdcdafasd
2 16 002agdcxxxxxx
3 23 001adfadfahoh
4 23 001fsdadggggg
```

得到：

```
1 16 001
2 16 002
3 23 001
4 23 002
```

```
1 awk '{print $1,substr($2,1,3)}'  
2 awk 'BEGIN{FIELDWIDTH="2 2:3"}{print $1,$2}' a.txt
```

## 行列转换

```
1 name age  
2 alice 21  
3 ryan 30
```

转换得到:

```
1 name alice ryan  
2 age 21 30
```

```
1 awk '  
2 {  
3     for(i=1;i<=NF;i++){  
4         if(!(i in arr)){  
5             arr[i]=$i  
6         } else {  
7             arr[i]=arr[i]" "$i  
8         }  
9     }  
10 }  
11 END{  
12     for(i=1;i<=NF;i++){  
13         print arr[i]  
14     }  
15 }  
16 ' a.txt
```

## 行列转换2

文件内容:

```
1 74683 1001  
2 74683 1002  
3 74683 1011  
4 74684 1000  
5 74684 1001  
6 74684 1002  
7 74685 1001  
8 74685 1011  
9 74686 1000  
10 ....  
11 100085 1000  
12 100085 1001
```

文件就两列，希望处理成



```

1 74683 1001 1002 1011
2 74684 1000 1001 1002
3 ...

```

就是只要第一列数字相同，就把他们的第二列放一行上，中间空格分开

```

1 {
2     if($1 in arr){
3         arr[$1] = arr[$1]" "$2
4     } else {
5         arr[$1] = $2
6     }
7
8 }
9
10 END{
11     for(i in arr){
12         printf "%s %s\n",i,arr[i]
13     }
14 }

```

## 筛选给定时间范围内的日志

grep/sed/awk用正则去筛选日志时，如果要精确到小时、分钟、秒，则非常难以实现。

但是awk提供了mktime()函数，它可以将时间转换成epoch时间值。

```

1 # 2019-11-10 03:42:40转换成epoch
2 $ awk 'BEGIN{print mktime("2019 11 10 03 42 40")}'
3 1573328560

```

借此，可以取得日志中的时间字符串部分，再将它们的年、月、日、时、分、秒都取出来，然后放入mktime()构建对应的epoch值。因为epoch值是数值，所以可以比较大小，从而决定时间的大小。

下面strptime1()实现的是将 2019-11-10T03:42:40+08:00 格式的字符串转换成epoch值，然后和which\_time比较大小即可筛选出精确到秒的日志。

```

1 BEGIN{
2     # 要筛选什么时间的日志，将其时间构建成epoch值
3     which_time = mktime("2019 11 10 03 42 40")
4 }
5
6 {
7     # 取出日志中的日期时间字符串部分
8     match($0,"^.*\\[(.*)\\].*",arr)
9
10    # 将日期时间字符串转换为epoch值
11    tmp_time = strptime1(arr[1])
12
13    # 通过比较epoch值来比较时间大小
14    if(tmp_time > which_time){print}
15 }

```

```

16
17 # 构建的时间字符串格式为: "2019-11-10T03:42:40+08:00"
18 function strtptime1(str ,arr,Y,M,D,H,m,S) {
19     patsplit(str,arr,"[0-9]{1,4}")
20     Y=arr[1]
21     M=arr[2]
22     D=arr[3]
23     H=arr[4]
24     m=arr[5]
25     S=arr[6]
26     return mktime(sprintf("%s %s %s %s %s %s",Y,M,D,H,m,S))
27 }

```

下面strptime2()实现的是将 10/Nov/2019:23:53:44+08:00 格式的字符串转换成epoch值, 然后和 which\_time比较大可即可筛选出精确到秒的日志。

```

1 BEGIN{
2     which_time = mktime("2019 11 10 03 42 40")
3 }
4
5 {
6     match($0,"^.*\\[(.*)\\].*",arr)
7
8     tmp_time = strtptime2(arr[1])
9
10    if(tmp_time > which_time){
11        print
12    }
13 }
14
15 # 构建的时间字符串格式为: "10/Nov/2019:23:53:44+08:00"
16 function strtptime2(str ,dt_str,arr,Y,M,D,H,m,S) {
17     dt_str = gensub("[/:+]", " ", "g",str)
18     # dt_sr = "10 Nov 2019 23 53 44 08 00"
19     split(dt_str,arr," ")
20     Y=arr[3]
21     M=mon_map(arr[2])
22     D=arr[1]
23     H=arr[4]
24     m=arr[5]
25     S=arr[6]
26     return mktime(sprintf("%s %s %s %s %s %s",Y,M,D,H,m,S))
27 }
28
29 function mon_map(str ,mons){
30     mons["Jan"]=1
31     mons["Feb"]=2
32     mons["Mar"]=3
33     mons["Apr"]=4
34     mons["May"]=5
35     mons["Jun"]=6
36     mons["Jul"]=7
37     mons["Aug"]=8
38     mons["Sep"]=9

```

```

39     mons["Oct"]=10
40     mons["Nov"]=11
41     mons["Dec"]=12
42     return mons[str]
43 }

```

## 去掉 **/\*\*/** 中间的注释

示例数据：

```

1  /*AAAAAAAAAA*/
2  1111
3  222
4
5  /*aaaaaaaa*/
6  32323
7  12341234
8  12134 /*bbbbbbbbbb*/ 132412
9
10 14534122
11 /*
12     cccccccccc
13 */
14 xxxxxx /*ddddddddddd
15     cccccccccc
16     eeeeeeee
17 */ yyyyyyyy
18 5642341

```

```

1  # 注释内的行
2  /\/*/{
3      # 同行有"*/"
4      if(/\*\/){
5          print gensub("(.*)/\*.*\*/(.*)", "\\1\\2", "g", $0)
6      } else {
7          # 同行没有"*/"
8
9          # 1.去掉*/行后的内容
10         print gensub("(.*)/\*.*", "\\1", "g", $0)
11
12         # 2.继续读取,直到出现*/,并去掉中间的所有数据
13         while( ( getline ) > 0 ){
14             # 出现了*/行
15             if(/\*\/){
16                 print gensub(".*\*/(.*)", "\\1", "g", $0)
17             }
18         }
19     }
20 }
21 # 非注释内容
22 !/\/*/{print}

```

## 前后段落关系判断

从如下类型的文件中，找出false段的前一段为i-order的段，同时输出这两段。

```
1 2019-09-12 07:16:27 [-][
2   'data' => [
3     'http://192.168.100.20:2800/api/payment/i-order',
4   ],
5 ]
6 2019-09-12 07:16:27 [-][
7   'data' => [
8     false,
9   ],
10 ]
11 2019-09-21 07:16:27 [-][
12   'data' => [
13     'http://192.168.100.20:2800/api/payment/i-order',
14   ],
15 ]
16 2019-09-21 07:16:27 [-][
17   'data' => [
18     'http://192.168.100.20:2800/api/payment/i-user',
19   ],
20 ]
21 2019-09-17 18:34:37 [-][
22   'data' => [
23     false,
24   ],
25 ]
```

```
1 BEGIN{
2   RS="]\n"
3   ORS=RS
4 }
5 {
6   if(/false/ && prev ~ /i-order/){
7     print tmp
8     print
9   }
10  tmp=$0
11 }
```

## 两个文件的处理

有两个文件file1和file2，这两个文件格式都是一样的。

需求：先把文件2的第五列删除，然后用文件2的第一列减去文件一的第一列，把所得结果对应的贴到原来第五列的位置，请问这个脚本该怎么编写？

```
1 file1:
```

```

2  50.481  64.634  40.573  1.00  0.00
3  51.877  65.004  40.226  1.00  0.00
4  52.258  64.681  39.113  1.00  0.00
5  52.418  65.846  40.925  1.00  0.00
6  49.515  65.641  40.554  1.00  0.00
7  49.802  66.666  40.358  1.00  0.00
8  48.176  65.344  40.766  1.00  0.00
9  47.428  66.127  40.732  1.00  0.00
10 51.087  62.165  40.940  1.00  0.00
11 52.289  62.334  40.897  1.00  0.00
12 file2:
13 48.420  62.001  41.252  1.00  0.00
14 45.555  61.598  41.361  1.00  0.00
15 45.815  61.402  40.325  1.00  0.00
16 44.873  60.641  42.111  1.00  0.00
17 44.617  59.688  41.648  1.00  0.00
18 44.500  60.911  43.433  1.00  0.00
19 43.691  59.887  44.228  1.00  0.00
20 43.980  58.629  43.859  1.00  0.00
21 42.372  60.069  44.032  1.00  0.00
22 43.914  59.977  45.551  1.00  0.00

```

```

1  # 方法一:
2
3  awk '{
4      f1 = $1
5      if( (getline <"file2") >= 0 ){
6          $5 = $1 - f1
7          print $0
8      }
9  }' file1
10
11 # 方法二:
12 awk '
13     NR==FNR{arr[FNR]=$1}
14     NR!=FNR{$5=$1-arr[FNR];print}
15 ' file1 file2

```