

19 文本处理工具awk编程

文本处理工具awk编程

awk 的工作流程：`awk '[BEGIN{动作}] 模式{动作} [END {动作}]' 文件...`

一、流程控制语句

流程控制语句是所有程序设计语言必不可少的部分，每一门语言都支持某些执行流程控制的语句。

awk 提供了完整的流程控制语句。

1.1 if条件语句

`if...else` 语句其语法格式为

```
1  if (条件表达式)
2      语句1
3  else
4      语句2
```

Shell | 复制代码

或：

```
1  if (条件表达式) 语句1; else 语句2
```

Shell | 复制代码

如果条件表达式的判断结构为真，则执行语句1，否则执行语句2。

格式中的 `语句1` 可以是**多个语句**，为了方便 `awk` 判断，可以将**多个语句用 `{}` 括起来**。

awk 分支结构允许嵌套，其格式为

```
1  if (条件表达式1)
2  {if (条件表达式2)
3      语句1
4  else
5      语句2
6  }
7  语句3
8  else {if (条件表达式3)
9      语句4
10 else
11     语句5
12 }
13 语句6
```

案例1：根据UID判断系统中用户类别

```
1 # 根据/etc/passwd第三列UID的大小判断是否为系统用户
2 # 1 ~ 200: 由Linux发行版自行建立的系统账号
3 # 201 ~ 999: 若用户有系统账号需求时, 就可以使用的账号UID
4 # 1000以上: 普通用户
5 [root@Shell ~]# awk -F: '{if($3<=200){name="system"} else{name="user"} print $1,name}' /etc/passwd
6 root system
7 bin system
8 daemon system
9 adm system
10 lp system
11 sync system
12 shutdown system
13 halt system
14 mail system
15 operator system
16 games system
17 ftp system
18 nobody system
19 systemd-network system
20 dbus system
21 polkitd user
22 sshd system
23 postfix system
24 chrony user
25
26 # 语句过长可用\续行
27 [root@Shell ~]# awk -F: '{if($3<=200){name="system"}else{name="user"} \
28 > print $1,name}' /etc/passwd
29 root system
30 bin system
31 daemon system
32 adm system
33 lp system
34 sync system
35 shutdown system
36 halt system
37 mail system
38 operator system
39 games system
40 ftp system
41 nobody system
42 systemd-network system
43 dbus system
44 polkitd user
```

```
45 sshd system
46 postfix system
47 postfix system
48 chrony user
49
49 # 模式部分可自动续行, 注意'位置
50 [root@Shell ~]# awk -F: '
51 > {
52 > if($3<=200)
53 > {name="system"}
54 > else
55 > {name="user"}
56 > print $1,name
57 > }' /etc/passwd
58 root system
59 bin system
60 daemon system
61 adm system
62 lp system
63 sync system
64 shutdown system
65 halt system
66 mail system
67 operator system
68 games system
69 ftp system
70 nobody system
71 systemd-network system
72 dbus system
73 polkitd user
74 sshd system
75 postfix system
76 chrony user
```

案例2：统计系统用户数

```

1 [root@Shell ~]# awk -F: '{if ($3>0 && $3<1000){++i}}{print i}' /etc/passwd
2
3 1
4 2
5 3
6 4
7 5
8 6
9 7
10 8
11 9
12 10
13 11
14 12
15 13
16 14
17 15
18 16
19 17
20 18
21 # END在所有行处理完才输出
22 [root@Shell ~]# awk -F: '{if ($3>0 && $3<1000){++i}} END {print i}' /etc/passwd
23 18

```

案例3：分别统计管理员个数，普通用户个数及系统用户个数

```

1 [root@Shell ~]# awk -F: '{if ($3==0){i++} else if ($3>999){j++} else {k++}}
  END {print "管理员"i,"普通用户"j,"系统用户"k}' /etc/passwd
2 管理员1 普通用户 系统用户18

```

1.2 while循环

`while` 循环语法结构为：

```

1 while 条件表达式
2 {语句}

```

案例4：while循环

```
▼ Shell 复制代码
1  # 注意，模式为BEGIN
2  [root@shell ~]# awk 'BEGIN {i=1; while(i<=10) {print i;i++}}'
3  1
4  2
5  3
6  4
7  5
8  6
9  7
10 8
11 9
12 10
```

案例5：使用while循环打印文件/etc/passwd中以root开头的行出每个字段

```
▼ Shell 复制代码
1 [root@shell ~]# awk -F: '/^root/ {i=1;while(i<=NF){print $i;i++}}' /etc/passwd
2 root
3 x
4 0
5 0
6 root
7 /root
8 /bin/bash
```

`awk` 除 `while` 循环结构外，还有 `do...while` 循环结构。它在代码块结尾处对条件求值，而不像标准 `while` 循环那样在开始处求值，其语法结构为：

```
▼ Shell 复制代码
1 do
2     语句
3 while (条件表达式)
```

与一般的 `while` 循环不同，由于在代码块之后对条件求值，`do...while` 循环结构永远都至少执行一次。

换句话说，当第一次遇到普通 `while` 循环时，如果条件为假，将永远不执行该循环。

案例6：使用do while循环打印文件/etc/passwd中以root开头的行出每个字段

▼ Shell 复制代码

```
1 [root@shell ~]# awk -F: '/^root/ {i=1;do {print $i;i++} while(i<=NF)}' /etc/passwd
2 root
3 x
4 0
5 0
6 root
7 /root
8 /bin/bash
```

1.3 for循环

`for` 循环中数组遍历的格式为：

▼ Shell 复制代码

```
1 for (变量 in 数组)
2 {语句}
```

`for` 循环中固定循环的格式为：

▼ Shell 复制代码

```
1 for (变量;条件;表达式)
2 {语句}
```

`for` 语句首先执行初始化语句，然后再检查条件。如果条件为真，则执行语句，然后执行递增或者递减操作。只要条件为真，循环就会一直执行。每次循环结束都会进行条件检查，若条件为假则结束循环。

案例7：使用for循环输出数字1至5

```
1 [root@shell ~]# awk 'BEGIN {for(i=1;i<=5;i++) {print i}}'
```

2	1
3	2
4	3
5	4
6	5

案例8：使用for循环打印文件/etc/passwd中以root开头的行出每个字段

```
1 [root@shell ~]# awk -F: '/^root/ {for(i=1;i<=NF;i++){print $i}} ' /etc/passwd
```

2	root
3	x
4	0
5	0
6	root
7	/root
8	/bin/bash

1.4 break命令、continue命令、exit命令

在 `awk` 的 `while` 语句、`do...while` 和 `if` 语句中可以使用 `break` 命令、`continue` 命令控制流程走向。`break` 中断当前正在执行的循环并跳到循环外执行下一条语句，`continue` 用于在循环体内部结束本次循环，从而直接进入下一次循环迭代，`exit` 用于结束脚本程序的执行。

案例9：当累加和大于50时，使用break结束循环


```
1 [root@shell ~]# awk 'BEGIN{ sum=0;for(i=0;i<20;++i){sum +=i;if(sum>50)break;else print"Sum=",sum}}'
```

Line	Sum
2	0
3	1
4	3
5	6
6	10
7	15
8	21
9	28
10	36
11	45

案例10：输出1到20之间的偶数

```
1 [root@shell ~]# awk 'BEGIN{for(i=1;i<=20;++i){if(i%2==0) print i;else continue}}'
```

Line	Output
2	2
3	4
4	6
5	8
6	10
7	12
8	14
9	16
10	18
11	20

案例11：当计算的和大于50时，结束awk程序

```

1 [root@shell ~]# awk 'BEGIN{ sum=0;for(i=0;i<20;++i){sum +=i;if(sum>50) exit(10);else print"Sum=",sum}}'
2 Sum= 0
3 Sum= 1
4 Sum= 3
5 Sum= 6
6 Sum= 10
7 Sum= 15
8 Sum= 21
9 Sum= 28
10 Sum= 36
11 Sum= 45
12 [root@shell ~]# echo $?
13 10

```

二、数组

`awk` 处理文本中数组是必不可少的，由于数组索引(下标)可以是**数字和字符串**，索引(下标)一般称作 `key`，并且与对应数组元素的值**关联**。因此，`awk` 中的**数组**称为**关联数组**(Associative Arrays)。

数组元素的 `key` 和值都放在 `awk` 内部程序的某一张表中，通常使用一定**散列算法**来存放，所以**数组元素并不是按照一定顺序来放**的。同理，也不是按照一定的顺序打印出来的，但可以使用管道来对所需的数据再次操作来达到效果。

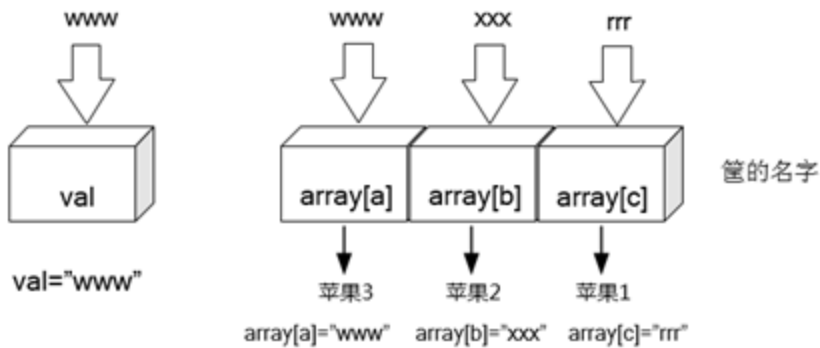
`awk` 中的数组**不必提前声明**，也**不必声明大小**，因为它在运行时可以自动地增加或减少。数组元素用 `0` 或**空字符串**来初始化，这根据上下文而定。

通常 `awk` 数组用来从记录中收集信息，可以用于计算总和、统计单词以及跟踪模板被匹配的次数等。

`awk` 数组的结构如下所示。

数组名 酒店名称	元素名 (房间号码)	值 (房间里面的人)
<code>arrayname</code>	<code>[string]</code>	<code>=value</code>

如上图可以发现，`awk` 数组就和酒店一样，数组的名称就像是酒店名称，数组元素名称就像酒店房间号码，每个数组元素里面的内容就像是酒店房间里的人。



定义数组的格式为：`数组名[下标]=元素值`，其中索引下标可以是数字，也可以是字符。

使用数组的格式为：`数组名[下标]`

输出数组元素的值：`print 数组名[下标]`

遍历数组循环结构：`for(变量名 in 数组名){print 数组名[变量名]}`

案例12：遍历数组，显示/etc/passwd的root账户

```
Shell | 复制代码
1 [root@shell ~]# awk -F: '{username[++i]=$1} END{print username[1]}' /etc/passwd
2 root
3 [root@shell ~]# awk -F: '{username[++i]=$1} END{print username[2]}' /etc/passwd
4 bin
```

案例12：遍历数组，显示/etc/passwd的所有账户

```
1 [root@shell ~]# awk -F: '{user[j++]=$1} END{for (i in user){print i,user  
[i]}}' /etc/passwd  
2 17 postfix  
3 4 lp  
4 5 sync  
5 6 shutdown  
6 7 halt  
7 8 mail  
8 9 operator  
9 10 games  
10 11 ftp  
11 12 nobody  
12 13 systemd-network  
13 0 root  
14 14 dbus  
15 1 bin  
16 15 polkitd  
17 2 daemon  
18 16 sshd  
19 3 adm
```

案例13：使用数组统计文件/etc/passwd中各种类型Shell的数量

```
1 [root@shell ~]# awk -F: '{shells[$NF]++} END{for (i in shells)print i,shell  
s[i]}' /etc/passwd  
2 /bin/sync 1  
3 /bin/bash 1  
4 /sbin/nologin 14  
5 /sbin/halt 1  
6 /sbin/shutdown 1
```

三、函数

定义、调用用户本身的函数是每个高级语言都具有的功能，`awk` 也不例外。

`awk` 函数分为**内置函数**和**自定义函数**。

3.1 内置函数

常用字符串内置函数如下表所示。

内置函数	作用
<code>gsub(x,y,z)</code>	在字串z中使用字串y替换与正则表达式x相匹配的所有字串，z默认为\$0。相当于sed中的s ///g
<code>sub(x,y,z)</code>	在字串z中使用字串y替换与正则表达式x相匹配的第一个字串，z默认为\$0。相当于sed中的 s///
<code>length(string)</code>	返回string参数指定的字符串的长度（字符形式）。如果未给出string参数，则返回整个记录的长度（\$0记录变量）。
<code>getline</code>	从输入中读取下一行内容。
<code>index(string1,string2)</code>	在由string1参数指定的字符串（其中有出现string2指定的参数）中，返回位置，从1开始编号。如果string2参数不在string1参数中出现，则返回0。
<code>substr(string,M,[N])</code>	返回具有N参数指定的字符数量字串。字串从string参数指定的字符串取得，其字符以M参数指定的位置开始。M参数指定为将string参数中的第一个字符作为编号1。如果未指定N参数，则字串的长度将是M参数指定的位置到string参数的末尾的长度。
<code>match(string,Ere)</code>	在string参数指定的字符串（Ere参数指定的扩展正则表达式出现在其中）中返回位置（字符形式），从1开始编号，或如果Ere参数不出现，则返回0（零）RSTART特殊变量设置为返回值，RLENGTH特殊变量设置为匹配的字符串的长度，或如果未找到任何匹配，则设置为-1（负一）。
<code>split(string,A,[Ere])</code>	将string参数指定的参数分割为数组元素A[1],A[2],...A[n]，并返回n变量的值。此分隔可以通过Ere参数指定的扩展正则表达式进行，或用当前字段分隔符（FS特殊变量）来进行（如果没有给出Ere参数）。除非上下文指明特定的元素还应具有一个数字值，否则A数组的元素用字符串值来创建。
<code>tolower(string)</code>	返回string参数指定的字符串，字符串中每个大写字符将更改为小写。

<code>sprintf(Format,Expr,Expr,...)</code>	根据Format参数指定的printf子例程格式字符串来格式化Expr参数指定的表达式并返回最后生成的字符串。
--	---

常用算术内置函数如下表所示。

内置函数	说明
<code>rand()</code>	产生0-1之间的浮点类型的随机数， <code>rand()</code> 产生随机数时需要设置一个参数，否则单独的 <code>rand()</code> 每次产生的随机数都是一样的。
<code>init(x)</code> <code>;</code>	返回 <code>x</code> 的整数部分的值。
<code>sqrt()</code>	返回 <code>x</code> 的平方根。
<code>srand()</code> <code>;</code>	建立 <code>rand()</code> 新的种子数，如果没有指定就用当天的时间。

案例14：统计用户名为4个字符的用户

```
1 # 传统方法
2 [root@shell ~]# awk -F: '$1~/^....$/ {count++;print $1} END{print count}'
   /etc/passwd
3 root
4 sync
5 halt
6 mail
7 dbus
8 sshd
9 6
10 # 使用awk内建函数length()计算并返回字符串的长度
11 [root@shell ~]# awk -F: 'length($1)==4 {count++;print $1} END{print count}' /etc/passwd
12 root
13 sync
14 halt
15 mail
16 dbus
17 sshd
18 6
```

3.2 自定义函数

函数定义包括函数名、函数参数、函数体。

`awk` 函数的定义方法如下：

```
1 function 函数名 (参数表) {
2     函数体
3 }
```

函数名是用户自定义函数的名称，函数名称以字母开头，后面可以是数字、字母或下画线的自由组合，`awk` 保留的关键字不能作为用户自定义函数的名称。

自定义函数可以接受参数，参数之间用逗号分隔，**参数不是必需的**，用户也可以定义没有任何输入参数的函数，函数体包含 `awk` 程序代码。

案例15：简易函数

```
1 [root@shell ~]# vi awk_test
2 11 22
3 [root@shell ~]# awk 'function test(a,b) { printf a "\t" b "\n"}
4 > { test($1,$2) } ' awk_test
5 11      22
```

当函数代码量比较大时，将函数写在外部文件中更合适。

案例16：返回两个数中的最大值和最小值

```
1 [root@shell ~]# vi functions.awk
2 function find_min(num1,num2)
3 {
4     if (num1<num2)
5         return num1
6     else
7         return num2
8 }
9
10 function find_max(num1,num2)
11 {
12     if (num1>num2)
13         return num1
14     else
15         return num2
16 }
17
18 function main(num1,num2)
19 {
20     result1=find_max(num1,num2)
21     result2=find_min(num1,num2)
22     print "最小值: ",result1,"最大值: ",result2
23 }
24
25 BEGIN{
26     main(10,20)
27 }
28 [root@shell ~]# awk -f functions.awk
29 最大值: 20 最小值: 10
```

案例17：自定义函数应用

上面的案例函数的参数是固定的，下面演示如何通过文件向函数传递参数。

```
▼ Shell | 复制代码

1 ▾ [root@shell ~]# vi functions.awk
2  function find_min(num1,num2)
3  {
4      if (num1<num2)
5          return num1
6      else
7          return num2
8  }
9
10 function find_max(num1,num2)
11 {
12     if (num1>num2)
13         return num1
14     else
15         return num2
16 }
17
18 function main(num1,num2)
19 {
20     result1=find_max(num1,num2)
21     result2=find_min(num1,num2)
22     print "最大值: ",result1,"最小值: ",result2
23 }
24
25 main($1,$2)
26
27 ▾ [root@shell ~]# awk -f functions.awk awk_test
28 最大值:  22 最小值:  11
29  # 修改样例文件中数据的分隔符
30 ▾ [root@shell ~]# vi awk_test
31 11:22
32 ▾ [root@shell ~]# awk -f functions.awk -F":" awk_test
33 最大值:  22 最小值:  11
34  # 增加样例文件中的记录
35 ▾ [root@shell ~]# vi awk_test
36 11:22
37 33:44
38 ▾ [root@shell ~]# awk -f functions.awk -F":" awk_test
39 最大值:  22 最小值:  11
40 最大值:  44 最小值:  33
```

注意：在无输入文件或标准输入时，函数调用部分要放在 `BEGIN` 模式中，为什么？

小结

- 流程控制：语法格式
- 数组：语法
- 函数：内置函数功能，自定义函数定义及调用

课程目标

- 知识目标：熟练掌握 `awk` 命令的基本语法。
- 技能目标：能够利用 `awk` 命令完成实战场景的处理。

课外拓展

- 进一步了解 `awk` 命令的应用场景。

参考资料

- `awk --help` 或 `man awk`
- 《Linux Shell核心编程指南》，丁明一，电子工业出版社