

# AWK 快速参考

模式、字段、数组、函数、文本处理

## 基础

### 运行 AWK

```
awk '{ print }' file.txt # print every line
awk '{ print $1 }' file.txt # print first field
awk -F: '{ print $1 }' /etc/passwd # custom delimiter
awk -f script.awk file.txt # run from file
cmd | awk '{ print $2 }' # pipe input
```

### 程序结构

<b>awk 'pattern { action }'</b>	基本格式——模式匹配时执行 action
<b>BEGIN { ... }</b>	在处理输入之前运行一次
<b>END { ... }</b>	处理完所有输入后运行一次
<b>无模式</b>	对每一行都执行 action
<b>无 action</b>	默认 action 为 { print }

## 模式与动作

### 模式类型

```
awk '/error/' file.txt # regex match
awk '$3 > 100' file.txt # comparison
awk 'NR >= 5 && NR <= 10' file.txt # line range
awk '/start/,/end/' file.txt # range pattern
```

### 模式参考

<b>/regex/</b>	将行与正则表达式匹配
<b>\$1 ~ /pat/</b>	字段匹配正则
<b>\$1 !~ /pat/</b>	字段不匹配正则
<b>expr1, expr2</b>	范围: 从第一个匹配到第二个
<b>expr1 &amp;&amp; expr2</b>	逻辑与
<b>expr1    expr2</b>	逻辑或
<b>!expr</b>	逻辑非

## 变量

### 内置变量

<b>NR</b>	当前记录 (行) 编号
<b>NF</b>	当前记录的字段数
<b>FS</b>	输入字段分隔符 (默认: 空白)
<b>OFS</b>	输出字段分隔符 (默认: 空格)
<b>RS</b>	输入记录分隔符 (默认: 换行)
<b>ORS</b>	输出记录分隔符 (默认: 换行)
<b>FILENAME</b>	当前输入文件名
<b>FNR</b>	当前文件中的记录编号

### 用户变量

```
awk '{ total += $1 } END { print total }' file.txt
awk -v threshold=50 '$1 > threshold' file.txt
awk 'BEGIN { count = 0 } /pat/ { count++ }
END { print count }' file.txt
```

## 字段

### 字段访问

<b>\$0</b>	整行内容
<b>\$1, \$2, ...</b>	第一、第二.....字段
<b>\$NF</b>	最后一个字段
<b>\$(NF-1)</b>	倒数第二个字段

### 字段分隔符

```
awk -F, '{ print $2 }' data.csv # comma
awk -F'\t' '{ print $1 }' data.tsv # tab
awk 'BEGIN { FS = "[::]" } { print $1 }' f # multi-char
awk 'BEGIN { OFS = "," } { print $1, $3 }' f # output sep
```

## 控制流

### 条件与循环

```
awk '{ if ($1 > 50) print "high"; else print "low" }' f
awk '{ for (i = 1; i <= NF; i++) print i }' f
awk '{ i = 1; while (i <= NF) { print i; i++ } }' f
awk '/skip/ { next } { print }' f # skip matching lines
```

### 控制语句

<b>if (cond) { ... } else { ... }</b>	条件判断
<b>for (i = 0; i &lt; n; i++) { ... }</b>	C 风格 for 循环
<b>for (key in array) { ... }</b>	遍历数组键
<b>while (cond) { ... }</b>	while 循环
<b>do { ... } while (cond)</b>	do-while 循环
<b>next</b>	跳到下一条输入记录
<b>exit</b>	停止处理, 执行 END 块

## 函数

### 用户定义函数

```
awk 'function max(a, b) {
    return (a > b) ? a : b
}
{ print max($1, $2) }' file.txt
```

### 数值函数

<b>int(x)</b>	截断为整数
<b>sqrt(x)</b>	平方根
<b>sin(x), cos(x)</b>	三角函数
<b>log(x), exp(x)</b>	自然对数与指数
<b>rand()</b>	0 到 1 之间的随机浮点数
<b>srand(seed)</b>	设置随机数种子

## 数组

### 关联数组

```
awk '{ count[$1]++ }
END { for (k in count) print k, count[k] }' f
awk '{ arr[NR] = $0 }
END { for (i = NR; i >= 1; i--) print arr[i] }' f
```

### 数组操作

<b>arr[key] = val</b>	设置元素
<b>arr[key]</b>	获取元素 (访问时自动创建)
<b>key in arr</b>	判断键是否存在
<b>delete arr[key]</b>	删除单个元素
<b>delete arr</b>	删除整个数组
<b>for (k in arr)</b>	遍历键 (无序)
<b>length(arr)</b>	元素数量 (gawk)

## 字符串函数

### 字符串参考

<b>length(s)</b>	字符串长度
<b>substr(s, start, len)</b>	子串 (从 1 开始计数)
<b>index(s, target)</b>	target 在 s 中的位置 (未找到返回 0)
<b>split(s, arr, sep)</b>	将字符串拆分到数组
<b>sub(pat, repl, s)</b>	替换第一个匹配
<b>gsub(pat, repl, s)</b>	替换所有匹配
<b>match(s, pat)</b>	正则匹配位置 (设置 RSTART、RLENGTH)
<b>tolower(s) / toupper(s)</b>	大小写转换
<b>sprintf(fmt, ...)</b>	格式化字符串 (类似 C printf)

## 字符串示例

```
awk '{ gsub(/old/, "new"); print }' f # sed-like replace
awk '{ print toupper($0) }' f # uppercase all
awk '{ print substr($0, 1, 40) }' f # truncate to 40
```

## I/O

### 输出

```
awk '{ print $1, $2 }' f # space-separated
awk '{ printf "%s,%d\n", $1, $2 }' f # formatted output
awk '{ print $1 > "out.txt" }' f # redirect to file
awk '{ print $1 >> "out.txt" }' f # append to file
```

### I/O 参考

<b>print</b>	带 ORS 输出 (默认换行)
<b>printf fmt, ...</b>	格式化输出 (不自动换行)
<b>print &gt; file</b>	重定向到文件
<b>print &gt;&gt; file</b>	追加到文件
<b>print   cmd</b>	管道输出到命令
<b>getline &lt; file</b>	从文件读取一行
<b>cmd   getline var</b>	将命令输出读入变量
<b>close(file)</b>	关闭文件或管道

## 常用模式

### 单行命令

```
awk '{ sum += $1 } END { print sum }' f # sum column
awk 'END { print NR }' f # count lines
awk '!seen[$0]++' f # remove dupes
awk 'NF' f # remove blanks
awk '{ print NF }' f # fields per line
```

### 实用配方

<b>CSV 转 TSV</b>	<code>awk -F, 'BEGIN{OFS="\t"} {\$1=\$1; print}'</code>
<b>第 2 列求和</b>	<code>awk '{ s += \$2 } END { print s }'</code>
<b>前 N 行</b>	<code>awk 'NR &lt;= 10' (类似 head)</code>
<b>频次统计</b>	<code>awk '{ c[\$1]++ } END { for (k in c) print k, c[k] }'</code>
<b>区间提取</b>	<code>awk '/BEGIN/,/END/'</code>
<b>打印第 N 列</b>	<code>awk '{ print \$N }'</code> (N 替换为列号)