

Linux 操作系统及应用

第六章 — 数据处理工具

唐晓晟 李亦农

`txs@bupt.edu.cn` `hoplee@bupt.edu.cn`

BEIJING UNIVERSITY OF POSTS AND TELECOMMUNICATIONS (BUPT)
SCHOOL OF INFORMATION AND COMMUNICATION ENGINEERING



内容简介

- ① 高级命令介绍
- ② 正则表达式 (Regular Expression)
- ③ `grep` 家族
- ④ 附录：特殊符号的名字



cut |

- 语法: `cut {SECTION SPECIFIERS} [OPTIONS] [FILES]`
- 说明: `cut`命令浏览每行输入, 将每行分成节, 然后显示指定的节。如果不指定输入, 则使用`STDIN`。
- 节描述符

`-bLIST` 显示由`LIST`指定位置上的字节

`-cLIST` 显示由`LIST`指定位置上的字符

`-fLIST` 显示由`LIST`指定位置上的字段, 字段定界符有行起始、行结束和制表符



cut ||

- 你可以用几种不同的方式来组成LIST参数:
 - 使用逗号表示单个不连续的小节: `-f4,7`
 - 使用范围表示连续的几个小节: `-b2-4`
 - 使用`n-`表示从`n`小节开始直到行结束的所有小节。
- 选项:

<code>-dDELIMITER</code>	指定节描述符 <code>-f</code> 所使用的字段定界符
<code>-n</code>	用于支持 Unicode 字符集
<code>-s</code>	和节描述符 <code>-f</code> 一起使用, 用于删除不包含字段分界符的行



cut III

- 范例:

```
1 [Apple] $ more test1.txt
2 test2alongsentence
3 this is test1
4 [Apple] $ cut -f2,3 -d" " test1.txt
5 test2alongsentence
6 is test1
7 [Apple] $ cut -f2 -d" " -s test1.txt
8 is
9 [Apple] $ cut -c1-6 test1.txt
10 test2a
11 this i
```



find |

- 语法: `find [PATHS] [EXPRESSION]`
- 说明: `find`命令检查目录树并计算所遇到的每一个表达式的值。
- 缺省路径是当前路径, 缺省表达式是`-print`
- `find`根据下列规则来分割路径和表达式, 在命令行上第一个`-()`, `!`之前的部份为路径, 之后的是表达式。
- 表达式分为选项子句、测试子句和动作子句三部分。各个子句之间通过运算符连接, 缺省运算符为`-and`。



find ||

- 运算符:

(expr)	圆括号用于改变优先级
! expr	取反运算
-not expr	
expr1 expr2	
expr1 -a expr2	
expr1 -and expr2	短路与运算
expr1 -o expr2	短路或运算
expr1 -or expr2	
expr1, expr2	顺序运算, 返回expr2的值



find III

- 选项影响整个find命令，因此它们总是返回真值。

-daystart	当天的起点为00:00:00而不是现在
-depth	子目录深度优先搜索
-help	显示简短帮助信息
-maxdepth n	最多搜索n层子目录，n=0表示测试和动作仅限于命令行列出的目录本身
-mindepth n	最少搜索n层子目录，n=1表示所有除命令行上的目录都要处理
-mount	只检查和指定目录在同一个文件系统下的文件



find IV

■ find命令的测试:

<code>-amin n</code>	文件在n分钟之前被访问过
<code>-anewer FILE</code>	文件的 LAT 在FILE的 LAT 之后
<code>-atime n</code>	文件在n天之前被访问过
<code>-cmin n</code>	文件状态在n分钟之前被修改过
<code>-cnewer FILE</code>	文件状态的 LMT 在FILE的 LMT 之后
<code>-ctime n</code>	文件状态在n天之前被修改过
<code>-empty</code>	文件为普通文件或目录并且为空
<code>-false</code>	永假
<code>-fstype TYPE</code>	文件位于指定类型的文件系统上。可以用- <code>printf</code> 动作的%F参数显示支持的文件系统 类型
<code>-gid n</code>	文件属于 gid 为n的组



find V

<code>-group name</code>	文件属于名字为 <code>name</code> 的组。允许使用数字形式的 <code>gid</code> 来代替 <code>name</code>
<code>-links n</code>	文件有 <code>n</code> 个链接
<code>-mmin n</code>	文件内容在 <code>n</code> 分钟之前被修改过
<code>-mtime n</code>	文件内容在 <code>n</code> 天之前被修改过
<code>-newer FILE</code>	文件的 LMT 在 <code>FILE</code> 的 LMT 之后
<code>-name PATTERN</code>	文件名与模式 <code>PATTERN</code> 相匹配
<code>-iname PATTERN</code>	文件名与模式 <code>PATTERN</code> 相匹配且忽略大小写
<code>-path PATTERN</code>	目录名与模式 <code>PATTERN</code> 相匹配
<code>-ipath PATTERN</code>	目录名与模式 <code>PATTERN</code> 相匹配且忽略大小写
<code>-perm MODE</code>	文件的权限与 <code>MODE</code> 匹配。 <code>MODE</code> 的形式为 <code>+ -nnn</code> ， <code>+</code> 为按位或运算， <code>-</code> 为按位与运算
<code>-regex PATTERN</code>	文件的全名与模式 <code>PATTERN</code> 相匹配
<code>-iregex PATTERN</code>	文件的全名与模式 <code>PATTERN</code> 相匹配且忽略大小写



find VI

<code>-size n[bckw]</code>	文件使用n个单元的空间。 b 为字节， c 为字， k 为千字节， w 为双字
<code>-type C</code>	文件的类型为 C 。 C 的值可为： b 块设备文件、
<code>-xtype C</code>	
<code>-true</code>	永真
<code>-uid n</code>	文件属于 uid 为n的用户
<code>-user NAME</code>	文件属于名字为 NAME 的用户。允许使用数字形式的 uid 来代替 NAME
<code>-pid n</code>	进程号为n的文件



find VII

- 在上述测试中，出现数字 n 的地方（时间、大小）可以使用 $+n$ 表示大于 n ， $-n$ 表示小于 n ，而 n 则表示正好等于 n
- **find**命令的动作：

<code>-exec CMD \;</code>	在当前目录下执行命令 <code>CMD</code> ，返
<code>-ok CMD</code>	
<code>-print</code>	输出文件全名至 STDOUT
<code>-print0</code>	在输出的文件全名后加一个空字
<code>-fprint0 FILE</code>	
<code>-fprint FILE</code>	输出文件全名至指定的文件
<code>-printf FORMAT</code>	按给定的格式输出。格式描述字
<code>-fprintf FILE FORMAT</code>	



find VIII

例 1 将当前目录及其子目录下所有扩展名是 c 的文件列出来。

```
1 [Apple] $ find . -name "*.c"
```

例 2 将当前目录及其下子目录中所有普通文件列出来

```
1 [Apple] $ find . -type f
```

例 3 将目前目录及其子目录下所有最近 20 分钟内更新过的文件列出

```
1 [Apple] $ find . -ctime -20
```

例 4 将当前目录及其子目录下所有扩展名不是 c 的文件列出来:

```
1 [Apple]$ find . -name "*.c" -o print
```



find IX

例 5 显示目录的深度及文件的目录名和文件名：

```
1 [Apple]$ find ch1 -printf "%d      %h/%f\n"
2 0/ch1
3 1      ch1/ch1.cpp
4 1      ch1/ch1.h
5 1      ch1/ch1_main.cpp
```

例 6 显示所有具有全局可执行权限的文件：

```
1 [Apple]$ find . -type f -perm -001 -print
2 ./test
3 ./filecounter
4 ./filesplitor
5 ./dtree
```



find X

例 7 使用 find 命令来生成命令行（备份）：

```
1 [Apple]$ dir2="backup"
2 [Apple]$ find dir1 -type d -printf \
3           "mkdir -p '${dir2}/${P}'\n" -o \
4           ! -type d -printf \
5           "mv '%p' '${dir2}/${p}'\n"| sh
```

%P和%p都是输出文件名，%p的输出是从命令行上指定的搜索目录开始，而%P则去掉命令行上指定的目录



find XI

```
1 [Apple]$ tree
2 .
3 |-- dir2
4 |   |-- one.txt
5 |   `-- two.txt
6 |-- one.txt
7 `-- two.txt
8 [Apple]$ find dir1 -printf "%p\n"
9 dir1
10 dir1/one.txt
11 dir1/two.txt
12 dir1/dir2
13 dir1/dir2/one.txt
14 dir1/dir2/two.txt
```


find XII

```
15 [Apple]$ find dir1 -printf "%P\n"  
16 one.txt  
17 two.txt  
18 dir2  
19 dir2/one.txt  
20 dir2/two.txt
```



info

- 语法: **info** [OPTIONS] [NODENAMES]
- 说明: **info**命令启动 GNU 扩展超文本文档系统。系统中每一个小而精的信息被称为一个节点。



join |

- 语法: `join [OPTIONS] FILE1 FILE2`
- 说明: 将两个文件中拥有相同字段的行合并起来。
- 选项:

<code>-aSIDE</code>	<code>SIDE</code> 为 1 或 2 表示 <code>FILE1</code> 或 <code>FILE2</code> 。将 <code>SIDE</code> 侧文件中无法配对的行输出
<code>-eSTRING</code>	用 <code>STRING</code> 来代表没有的输入字段
<code>-i</code>	字段比较时忽略大小写
<code>-oFIELDLIST</code>	只显示出现在 <code>FIELDLIST</code> 中的字段。形式为用逗号分隔的 <code>SIDE.n</code>
<code>-tCHAR</code>	使用指定的字符作为输入和输出字段分隔符
<code>-vSIDE</code>	仅 <code>SIDE</code> 侧文件中无法配对的行输出
<code>-1 FIELD</code>	根据 <code>FILE1</code> 的指定字段来合并行
<code>-2 FIELD</code>	根据 <code>FILE2</code> 的指定字段来合并行



join II

- 范例：先建立三个示例文本文件：

```
1  jnames          jsalaries      jhiredt
2  -----
3  001 Barry       001 135000      001 04/12/1989
4  002 Mario       002 80000       002 06/23/1994
5  003 John        003 45000       003 12/12/1997
6  004 Harry       004 75000       004 11/17/1997
7  005 Dorothy     005 90000       005 02/06/1994
8  006 Jackie      007 68000       006 07/20/1995
9
10 [Apple]$ join jnames jsalaries
11 001 Barry 135000
12 002 Mario 80000
```

join III

```
13 003 John 45000
14 004 Harry 75000
15 005 Dorothy 90000
16 [Apple]$ join -a1 jnames jsalaries
17 001 Barry 135000
18 002 Mario 80000
19 003 John 45000
20 004 Harry 75000
21 005 Dorothy 90000
22 006 Jackie
23 [Apple]$ join -a2 jnames jsalaries
24 001 Barry 135000
25 002 Mario 80000
26 003 John 45000
27 004 Harry 75000
```

join IV

```
28 005 Dorothy 90000
29 007 68000
30 [Apple]$ join -o1.2,2.2 jnames jsalaries
31 Barry 135000
32 Mario 80000
33 John 45000
34 Harry 75000
35 Dorothy 90000
36 [Apple]$ join -v1 jnames jsalaries
37 006 Jackie
38 [Apple]$ join -v2 jnames jsalaries
39 007 68000
40 [Apple]$ join -a1 jnames jsalaries | \
41     join - jhiredt
42 001 Barry 135000 04/12/1989
```

join V

```
43 002 Mario 80000 06/23/1994
44 003 John 45000 12/12/1997
45 004 Harry 75000 11/17/1997
46 005 Dorothy 90000 02/06/1994
47 006 Jackie 07/20/1995
```



paste |

- 语法: `paste [OPTIONS] [FILES]`
- 说明: 将多个文件按行合并, 在中间加上一个TAB, 然后输出到标准输出上。
- 选项:

<code>-s</code>	依次先将第一个文件的所有行链接到一起, 然后输出一个回车再处理第二个文件
<code>-d LIST</code>	依次使用LIST中的字符取代TAB作为输出分隔符
<code>-f</code>	遇到第一个EOF 就终止

- 范例:



paste II

```
1 $ cat num2          $ cat let3
2 1                    a
3 2                    b
4                      c
5 $ paste num2 let3    $ paste -s num2 let3
6 1      a             1      2
7 2      b             a      b      c
8                      c
9 $ paste -d '%_' num2 let3 num2
10 1%a_1
11 2%b_2
12 %c_
```



nl

- 语法: `nl [OPTIONS] [FILES]`
- 说明: 加上行号后输出文件



pstree |

- 语法: `pstree [OPTIONS] [PID|USER]`
- 说明: 以树形格式输出进程名, 缺省时以`init`进程为根节点, 可以通过指定`PID`或`USER`来改变根节点。
- 选项:

-
- `-a` 显示各进程的命令行。换出的进程显示在括号种
 - `-c` 取消压缩等价的进程
 - `-h` 高亮当前进程及其子进程
 - `-l` 使用加长行显示, 避免被屏幕宽度所限
 - `-n` 根据进程标识符来排序, 缺省时按名称排序
 - `-p` 显示进程标识符
-



pstree II

- 范例:

```
1 [Apple]$ pstree -h
2 init--+-adsl-connect---pppd---pppoe
3      |-atd
4      .....
5      |-socks5---2*[socks5]
6      |-sshd--+-sshd---bash---ssh
7      |      `--sshd---bash---pstree
8      |-syslogd
9      `--xinetd
```



sort |

- 语法: `sort [OPTIONS] [FILES]`
- 说明: 将输入文件按行排序、将多个有序文件合并或验证文件是否是有序的。缺省时`sort`将单个空格作为字段定界符。然后根据最左边的字段对各行排序, 如果最左边的字段相同, 则检查下一个字段。
- 当给大文件排序时, 应保证有足够的临时磁盘空间。



sort II

- 选项:

<code>-b</code>	忽略每行开头的空格
<code>-d</code>	以电话簿顺序排序。忽略所有字母、数字和空格之外的字符
<code>-i</code>	忽略 ASCII(32-126 含) 字符集之外的所有字符
<code>-kn[,m]</code>	指定从 <code>n</code> 到 <code>m</code> 之前的字段为排序字段
<code>-m</code>	合并两个或更多的有序文件
<code>-n</code>	按数字排序
<code>-o OUTPUT</code>	将排序结果写入 <code>OUTPUT</code> 文件中
<code>-r</code>	反序
<code>-tC</code>	指定字段分隔符为 <code>C</code>
<code>-u</code>	相同的行只输出一次



sort III

■ 范例:

```
1 [Apple]$ cat date
2 1 AC BB CC
3 2 AB CC DD
4 3 CA BB CC
5 5 BE DD EE
6 4 BD AA AA
7 [Apple]$ sort date           [Apple]$ sort -k2 date
8 1 AC BB CC                  2 AB CC DD
9 2 AB CC DD                  1 AC BB CC
10 3 CA BB CC                 4 BD AA AA
11 4 BD AA AA                 5 BE DD EE
12 5 BE DD EE                 3 CA BB CC
13
```



sort IV

```
14 [Apple]$ sort -k2.3 date1 [Apple]$ sort -k2.3r date1
15 3 CA BB CC                    5 BE DD EE
16 2 AB CC DD                    4 BD AA AA
17 1 AC BB CC                    1 AC BB CC
18 4 BD AA AA                    2 AB CC DD
19 5 BE DD EE                    3 CA BB CC
20 [Apple]$ sort -n date1
21 1 AC BB CC
22 2 AB CC DD
23 3 CA BB CC
24 4 BD AA AA
25 5 BE DD EE
```



split

- 语法: `split [OPTIONS] [INPUT [PREFIX]]`
- 说明: `split`命令以大小作为划分的标准将文件分割成几个部分，分割后的文件名为 `PREFIXaa`、`PREFIXab`、...，缺省的前缀为 `x`，缺省的后缀长度为 2。
- 选项:

`-aN` 指定后缀长度为`N`

`-lNUM` 将输出文件的大小限制为`NUM`行

`-bNUM[b|k|m]` 将输出文件的大小限制为`NUM`字节，后续的`b`表示 512 字节、`k`为千字节、`m`为兆字节

`-cNUM[b|k|m]` 以整行为单位输出文件并且将输出文件的大小限制为`NUM`字节/512 字节/千字节/兆字节



strings

- 语法: **strings** [OPTIONS] [FILES]
- 说明: 此命令用于显示二进制文件中的可打印字符串。缺省时, **strings** 在文件的初始化小节和已装入小节中寻找所有长度大于 4 个字符的可打印字符序列。
- 选项:

-a	扫描整个文件查找字符串
-f	将文件名加到每个找到的字符串之前
-n LENGTH	指定可打印字符序列的最小长度
-o [o x d]	分别用八进制、十六进制和十进制输出每个字符串的位置



tee |

- 语法: `tee [OPTIONS] [FILES]`
- 说明: **T**型管道: 从**STDIN**中读取并写到 **STDOUT**和指定的文件中。
- 选项:

<code>-a</code>	追加写入文件
<code>-i</code>	忽略中断信号

- 范例:



tee ||

```
1 [Apple]$ ls -l | tee ls.dat | wc -l
2      12
3 [Apple]$ more ls.dat
4 total 52
5 drwxrwxr-x  2 hop hop 4096 Oct 16 21:43 ch1
6 -rw-rw-r--  1 hop hop   55 Nov  1 13:27 date1
7 drwxrwxr-x  3 hop hop 4096 Nov  1 11:33 dir1
8 -rwxrwxr-x  1 hop hop  195 Oct 31 22:56 dtree
9 -rwxrwxr-x  1 hop hop 3095 Oct 31 13:30 filecounter
10 -rwxrwxr-x  1 hop hop 1049 Oct 31 23:03 filesplitor
11 -rw-rw-r--  1 hop hop   90 Nov  1 12:15 jhiredt
12 -rw-rw-r--  1 hop hop   62 Nov  1 12:14 jnames
13 -rw-rw-r--  1 hop hop   61 Nov  1 12:17 jsalaries
14 -rwxrwxr-x  1 hop hop 1654 Oct 31 12:26 test
15 -rw-rw-r--  1 hop hop   33 Nov  1 09:38 test1.txt
```



tr |

- 语法: `tr [OPTIONS] CHARSET1 [CHARSET2]`
- 说明: 对从`STDIN`中读入的字符进行变换、压缩或删除, 然后将结果送至`STDOUT`, `CHARSET1`和`CHARSET2`指定了哪些字符被转换以及它们被转换成什么字符。
- 当`CHARSET1`和`CHARSET2`都被指定并且没有使用`-d`选项时, `tr`执行转换, 每一个与`CHARSET1`中的字符匹配的输入字符都被转换为`CHARSET2`中相应的字符。
- `tr`也可用于将重复出现的字符紧缩为一个字符。



tr ||

- 选项:

-
- c 将所有CHARSET1之外的字符转换为CHARSET2中的字符
 - d 删除CHARSET1中的字符，不作转换
 - s 将连续多个相同的字符紧缩为一个
 - t 截断CHARSET1使其与CHARSET2包含同样数目的字符
-



tr |||

- 转换时若CHARSET2比CHARSET1长，则忽略多余的字符；反之，则不断重复CHARSET2的最后一个字符直到二者一样长。
- [c*n]表示字符c重复n次，[c*]表示将字符c重复必要的次数以使CHARSET1和CHARSET2长度相等。
- 在CHARSET1和CHARSET2中可使用字符类。
- 范例：



tr IV

```
1 [Apple]$ cat tr.dat
2 ABC DEF GHI
3 ABC DEF GHI
4 AAA BBB CCC
5 [Apple]$ tr "A" "Z" < tr.dat
6 ZBC DEF GHI
7 ZBC DEF GHI
8 ZZZ BBB CCC
9 [Apple]$ tr "ABC" "[Z*]" < tr.dat
10 ZZZ DEF GHI
11 ZZZ DEF GHI
12 ZZZ ZZZ ZZZ
13 [Apple]$ tr "ABCD" "YYZZ" < tr.dat
14 YYZ ZEF GHI
15 YYZ ZEF GHI
16 YYY YYY ZZZ
17 [Apple]$ tr -s "BC" < tr.dat
```



tr V

```
18 ABC DEF GHI
19 ABC DEF GHI
20 AAA B C
21 [Apple]$ tr -s "BC" "ZZ" < tr.dat
22 AZ DEF GHI
23 AZ DEF GHI
24 AAA Z Z
25 [Apple]$ tr -d "\n" < tr.dat ; echo
26 ABC DEF GHIABC DEF GHIAAA BBB CCC
27 [Apple]$ tr -c "A" "Z" < tr.dat
28 AZZZZZZZZZZZAZZZZZZZZZZZAAAZZZZZZZZZ[Apple]$
29 [Apple]$ tr -c "A \n" "Z" < tr.dat
30 AZZ ZZZ ZZZ
31 AZZ ZZZ ZZZ
32 AAA ZZZ ZZZ
```



tty

- 语法: `tty [OPTIONS]`
- 说明: 输出连接到STDIN的终端的文件名。
- 选项:

`-s` 不产生任何输出，只返回一个退出值



uniq |

- 语法: `uniq [OPTIONS] [INPUT [OUTPUT]]`
- 说明: 从已排序的文件中删除重复的行。
- 选项:

<code>-c</code>	在输出每行之前显示该行重复的次数
<code>-d</code>	仅显示重复的行
<code>-fN</code>	在检查行的唯一性前跳过N个字段
<code>-sN</code>	在检查行的唯一性前跳过N个字符
<code>-u</code>	仅显示不重复的行



uniq II

■ 范例:

```
1 [Apple]$ cat uniq.dat      [Apple]$ uniq uniq.dat
2 01 12345 67890            01 12345 67890
3 02 12345 67890            02 12345 67890
4 03 12345 22222            03 12345 22222
5 03 12345 22222            04 11111 22222
6 04 11111 22222
7 [Apple]$ uniq -d uniq.dat  [Apple]$ uniq -f1 -w5 uniq.dat
8 03 12345 22222            01 12345 67890
9                             04 11111 22222
10 [Apple]$ uniq -u uniq.dat [Apple]$ uniq -f1 uniq.dat
11 01 12345 67890            01 12345 67890
12 02 12345 67890            03 12345 22222
13 04 11111 22222            04 11111 22222
```



xargs |

- 语法: `xargs [-0prt看] [-e[eof-str]] [-i[replace-str]] [-l[max-lines]] [-n max-args] [-s max-chars] [-P max-procs] [--null] [--verbose] [--exit] [--max-procs=max-procs] [--version] [--help] [COMMAND [INITIAL-ARGUMENTS]]`
- 说明: 建立并执行命令行。
- `xargs`从STDIN读入被空格和新行符分隔的参数, 然后以 INITIAL-ARGUMENTS后跟随读入的参数来执行给定的命令 (缺省时执行 `/bin/echo`)



xargs II

- 选项:

<code>-0</code>	输入文件名以空字符结束。取消空格、引号和反斜线的特殊含义
-----------------	------------------------------

`--null`

<code>-e[<i>eof-str</i>]</code>	指定文件结束字符串。缺省为下划线
---------------------------------	------------------

<code>-i[<i>replace-str</i>]</code>	用从 STDIN 读入的名字替换指定的字符串
-------------------------------------	-------------------------------

<code>-l[<i>max-lines</i>]</code>	每个命令行最多使用 max-lines 个非空的输入行，缺省为 1
-----------------------------------	--

<code>-n <i>max-args</i></code>	每个命令行最多使用 max-args 个参数
---------------------------------	-------------------------------

<code>-p</code>	提示用户是否执行每个命令行并读入下一输入行
-----------------	-----------------------

<code>-r</code>	如果 STDIN 不包含任何非空字符，则不执行命令
-----------------	----------------------------------



xargs III

`-s max-chars`

每个命令行最多使用`max-chars`个字符，包括命令本身、初始参数和标记结束的空字符，缺省上限为 20k 个字符

`-t`

`--verbose`

在命令执行前在`STDOUT`上输出整个命令行

`-x`

当`-l`、`-n`和`-s`选项指定的值被超出时终止本

`--exit`

`-P max-procs`

同时最多运行`max-procs`个进程，缺省值为 1。若`max-procs`被指定为 0，则进程数无限制



示例 I

统计单词出现频率

对于给定的英文文本文件，统计各个单词出现的频率，并且显示频率最高的 n 个单词及其出现次数。

■ 算法步骤：

- ① 将所有大小写英文字母之外的字符转化为一个回车并将多个回车压缩为一个
- ② 将所有大写字母转换为小写
- ③ 排序，结果是按字典顺序排好的单词列表，相同的单词挨在一起
- ④ 在每个单词之前显示此单词出现的次数
- ⑤ 按照出现的次数降序排序
- ⑥ 取前 n 行输出
- ⑦ 输出时加上序号



示例 II

- 脚本代码:

```
1  #!/bin/bash
2  # Author   : Hop Lee
3  # Date     : 2002.10
4  # Purpose  : Output the top n words in a given
5  #             English text file
6  #
7  if [ $# -ne 2 -a $# -ne 1 ]; then
8      echo "usage: `basename $0' [n] Input_File"
9      echo
10     exit
11 fi
12
13 if [ $# -eq 1 ]
14 then
15     I_TOP=10
```

示例 III

```
16     I_FILE=$1
17 fi
18
19 if [ $# -eq 2 ]
20 then
21     I_TOP=$1
22     I_FILE=$2
23 fi
24
25 tr -sc "[A-Z][a-z]" "[\012*]" < $I_FILE | \
26 tr "[A-Z]" "[a-z]" | \
27 sort | \
28 uniq -c | \
29 sort -k1 -n -r | \
30 head -$I_TOP | nl
```



正则表达式语法 I

- 正则表达式是一种记法，她可以让你搜索符合特定规范的文本。
- 我们可以通过构造单个正则表达式来选择、匹配多个字符串。
- 大量的 UNIX 命令使用了正则表达式，其中比较著名的有 **grep** 家族、**sed**、文本处理工具 (**awk**, **Perl**, **Python**, **Ruby**, **Tcl** 等)、文件查看工具 (**more**, **page**, **less** 等) 以及文本编辑器 (**ed**, **vi/vim**, **emacs**, **jed**, **jove** 等)。
- 正则表达式分为两大类: BRE (Basic Regular Expressions) 和 ERE (Extended Regular Expressions).



正则表达式语法 I

<code>x?</code>	0 或 1 个 <code>x</code> 字符
<code>x*</code>	0 或 0 个以上 <code>x</code> 字符
<code>.*</code>	0 或 0 个以上任意字符
<code>x+</code>	1 或 1 个以上 <code>x</code> 字符
<code>.+</code>	1 或 1 个以上任意字符
<code>{m}</code>	<code>m</code> 个字符
<code>{m,n}</code>	<code>m</code> 个以上、 <code>n</code> 个以下个数的字符
<code>{m,}</code>	<code>m</code> 个以上个数的字符
<code>[]</code>	范围 (若减号在最后, 则失去表示区间的意义)
<code>[~]</code>	范围的补集, 若 <code>~</code> 不在开头, 则失去其表示补集的特殊意义
<code>^</code>	字符串开头的字符
<code>\$</code>	字符串末尾的字符



正则表达式语法 II

<code>\d</code>	等价于 <code>[0-9]</code>
<code>\d+</code>	等价于 <code>[0-9] +</code>
<code>\D</code>	等价于 <code>[^0-9]</code>
<code>\D+</code>	等价于 <code>[^0-9] +</code>
<code>\w</code>	等价于 <code>[a-zA-Z0-9]</code>
<code>\w+</code>	等价于 <code>[a-zA-Z0-9] +</code>
<code>\W</code>	等价于 <code>[^a-zA-Z0-9]</code>
<code>\W+</code>	等价于 <code>[^a-zA-Z0-9] +</code>
<code>\s</code>	等价于 <code>[\n\t\r\f]</code>
<code>\s+</code>	等价于 <code>[\n\t\r\f] +</code>
<code>\S</code>	等价于 <code>[^\n\t\r\f]</code>
<code>\S+</code>	等价于 <code>[^\n\t\r\f] +</code>
<code>\b</code>	寻找不以数字和字母为边界的串
<code>\B</code>	寻找以数字和字母为边界的串
<code>a b c</code>	含有 <code>a</code> 或 <code>b</code> 或 <code>c</code> 的串



正则表达式语法 III

abc 含有 abc 的串

\(pattern\) 依次将找到的串存入标记寄存器 1,2,...,9

/pattern/i 不考虑大小写

\ 转义



正则表达式语法 IV

- **非特殊字符**与其自身匹配。
- 在**特殊字符**前，缀以**反斜线**，则丧失字符的特殊含义，与其自身匹配。
- **圆点、星号、反斜线**在**方括号**内时，代表它们自己。
- **\$**只有出现在正则表达式最尾部时才有特殊意义，否则与其自身匹配；**^**只有出现在正则表达式最左边时才有特殊意义，否则与其自身匹配。



正则表达式运算符

- `=~`: `$string=~/pattern/` 表示对字符串 `$string` 用模式 `/pattern/` 进行匹配运算
- `!~`: 不相匹配运算符
- 单目运算符`*`、`+`和`?`优先级最高，其次是串接符，最后是`|`。可以用圆括号改变运算的优先级。
- 所有运算都是左结合的。



常用的函数

<code>tr/srchlist/repllist/</code>	把srchlist转换为repllist
<code>s/pattern/replace/eg</code>	用replace替换pattern



标记正则表达式

- **标记正则表达式**是用\ (和\) 括起来的部分。
- 匹配时从左到右扫描并依次将匹配的结果存入19 号寄存器中，可以使用\ n 引用寄存器的内容。
- 例：'\([a-z]\)\([a-z]\)[0-9]\2\1'，这个标记正则表达式将匹配所有由 5 个字符组成、并且中间一位是数字的回文字符串。



正则表达式示例 I

- `/\b([a-z]+) \1\b/gi` 一个单词的连续出现 (以空格分隔);
- `/(\w+):\\\/([^\:]+)(:\d*)?([^\#]*)/` 将一个 URL 解析为协议、域、端口及相对路径;
- `/^(?:Chapter|Section)[1-9][0-9]{0,1}$` /定位章节的位置;
- `/[-a-z]/` A 至 z 共 26 个字母再加一个 - 号;
- `/ter\b/` 可匹配 chapter, 而不能匹配 terminal;
- `/\Bapt/` 可匹配 chapter, 而不能匹配 aptitude;
- `/Windows(95|98|NT)/` 可匹配 Windows95 或 Windows98 或 WindowsNT, 当找到一个匹配后, 从 Windows 后面开始进行下一次的检索匹配;
- `^(?=.*\d)(?=.*[a-z])(?=.*[A-Z]).{8,10}$` 验证强密码。字符数必须在 8 至 10 的范围内。必须包含大小写字母和数字的组合, 不能使用特殊字符;



正则表达式示例 II

- `\d+` 非负整数（正整数和 0）；
- `((-\d+)|(0+))` 非正整数（负整数和 0）；
- `[0-9]*[1-9][0-9]*` 正整数；
- `-[0-9]*[1-9][0-9]*` 负整数；
- `-?\d+` 整数；
- `\d+(\.\d+)?` 非负浮点数（正浮点数和 0）；
- `((-\d+(\.\d+)?)|(0+(\.\d+)?))` 非正浮点数（负浮点数和 0）；
- `(([0-9]+\.[0-9]*[1-9][0-9]*)|
([0-9]*[1-9][0-9]*\.[0-9]+)|([0-9]*[1-9][0-9]*))` 正浮点
数；
- `(-(([0-9]+\.[0-9]*[1-9][0-9]*)|
([0-9]*[1-9][0-9]*\.[0-9]+)|([0-9]*[1-9][0-9]*)))` 负浮点
数；



正则表达式示例 III

- `(-?\d+)(\.\d+)?` 浮点数;
- `[a-zA-Z][.?!]` 表示任意后面跟随句号、问号或感叹号的小写或大写字母;
- `[(){}]` 右方括号、右圆括号或右花括号;
- `[01][0-9][-/][0-3][0-9][-/][0-9][0-9]` 形如 MM-DD-YY 或 MM/DD/YY 的日期表示;
- `[^aeiou]` 任意非元音字母字符;
- `80[234]?86` 匹配 “8086”, “80286”, “80386” 或 “80486”
- `(^\. . .)` 表示 “行首为一个句点, 随后是两个字符, 然后跟一个空格的字符串”;
- `[0-9]\{3\}-[0-9]\{4\}` 北美地区电话号码;
- `compan(y|ies)` 匹配 “company” 的单数或复数形式;
- `[\u4e00-\u9fa5]` 汉字;



正则表达式示例 IV

- `[^\x00-\xff]` 任意双字节字符;
- `<[^>]*>` 任意的 HTML 标签;
- `<\(.*\)>.*</\1>|<(.*)/>` 更准确地匹配 HTML 标签;
- `\w+([-+.] \w+)*@\w+([-.] \w+)*\.\w+([-.] \w+)*` 匹配 Email 地址的正则表达式;
- `([\w-\.]+)@((([0-9]{1,3}\. [0-9]{1,3}\. [0-9]{1,3}\.)|(([\w-]+\.)+))([a-zA-Z]{2,4}|[0-9]{1,3}))(\?)` 更准确地匹配 Email 地址;
- `http://([\w-]+\.)+[\w-]+(/[\w- ./?%&=]*)?` 匹配 URL;
- `^(http|https|ftp)\://[a-zA-Z0-9\-\._]+\.[a-zA-Z]{2,3}(:[a-zA-Z0-9]*)?/?([a-zA-Z0-9\-\._]? \, \\' /\\ \+ & % \ $ # \ =) * $` 更准确地匹配 URL



grep家族概述

- **grep**命令可以在一个文件中搜索指定的字符串。其基本语法为：
grep [OPTIONS] [-e] 'PATTERN' [FILENAMES]
- **FILENAMES**是一系列用空格格开的文件名，输出结果时会在每个匹配之前显示文件名
- 如果没有**FILENAMES**参数，则**grep**命令将扫描标准输入
- 缺省情况下**grep**将输出匹配指定模式的每一行
- 在**PATTERN**中如果含有元字符的话必须将其转义，使用' 或 \ 或"
- **fgrep**命令是**grep**的一个简化版本，仅用于查找固定的字符串而不允许使用正则表达式。其优点为速度快，且可以指定任意个数的查找字符串。
- **egrep**命令具有**fgrep**的一些特征，同时也能使用正则表达式、增加了一些查找技巧，但是不支持**grep**和**ed**命令中的标记表达式 (Tag Expression)。



grep 家族选项总结 I

仅适用于 **fgrep** :

-x	严格匹配整行
-----------	--------

仅适用于 **grep** :

-s	隐藏错误信息
-----------	--------

-w	仅匹配整个字
-----------	--------

仅适用于 **egrep** 和 **fgrep** :

-f	匹配指定文件中的字符串
-----------	-------------



grep家族选项总结 II

适用于grep, fgrep和egrep :

- h 隐藏文件名的显示
- b 显示匹配行的磁盘块号
- l 仅显示发现匹配的文件名
- c 统计和显示匹配的数量
- n 显示每个匹配行的行号
- v 显示非匹配的行
- I 忽略大小写
- e 匹配后面的以“-”开头的表达式



grep家族表达式总结 I

仅适用于grep和egrep：

\m	转义元字符m
^	行首
\$	行尾
.	任意单个字符
[xy^\$z]	x,y,^,\$,z中的任意一个字符
[^xy^\$z]	上式的否定
[a-z]	指定的范围内的任一字符
[^a-z]	上式的否定
r*	零或多个r表达式
rs	匹配r和后续的s（连接）



grep 家族表达式总结 II

仅适用于 **grep** :

<code>\(r\)</code>	匹配标记正则表达式
--------------------	-----------

<code>\n</code>	设置为匹配第 n 个标记表达式 (1~9)
-----------------	------------------------------

仅适用于 **egrep** :

<code>r+</code>	匹配一个或多个 r
-----------------	------------------

<code>r?</code>	匹配零个或一个 r
-----------------	------------------

<code>r s</code>	匹配 r 或 s
------------------	------------------------

<code>(r s)t</code>	匹配 rt 或 st
---------------------	--------------------------

<code>(r s)*</code>	匹配零个或多个 r s
---------------------	--------------------



其它grep相关命令

agrep 可进行模糊搜索

ngrep 用于搜索网络层数据

pdfgrep 用于搜索pdf文件中的模式

ack 是**grep**的一个可选替换品，可以对匹配字符串高亮显示，默认递归搜索且能自动忽略无关的文件

sgrep 用于搜索结构化的模式，例如 SGML, C, TeX 等类型的源文件

xgrep 专门针对 XML 文件的搜索工具



grep 示例 I

- 示例数据文件:

```
1 [Apple]$ cat myfile
2 blueflybird is me.
3 BLUEFLYBIRD IS ME.
4 Yes,I'm blueflybird.
5 of course,Blueflybird is me.
6 Of Course,BLUEFLYBIRD IS ME.
7 [BLUEBLUE BLUE by BLUE blue bird bird ]
```



grep 示例 II

① 分支、块、原子

```
1 [Apple]$ egrep "^blueflybird|BLUEFLYBIRD.*E\.$" \
2     myfile
3 blueflybird is me.
4 BLUEFLYBIRD IS ME.
5 Of Course,BLUEFLYBIRD IS ME.
```

- 上面的命令包括两个分支，用|分开，意思是查找myfile文件中，以小写blueflybird开头的行，或者是含有大写BLUEFLYBIRD单词并且以E.结尾的行。



grep 示例 III

② 匹配行末\$

- 如果要查找myfile文件中以e.结尾的行（注意加了转义符\）:

```
1 [Apple]$ egrep "e\.$" myfile
2 blueflybird is me.
3 of course,Blueflybird is me.
```

- 或者匹配以bird.结尾的行（注意\$号放在命令后面）:

```
1 [Apple]$ egrep "bird\.$" myfile
2 Yes,I'm blueflybird.
```



grep 示例 IV

③ 匹配行首[^]

- 如果要查找myfile文件中以小写b开头的行:

```
1 [Apple]$ egrep "^b" myfile
2 blueflybird is me.
```

- 如想查找以小写b或者大写B开始的行:

```
1 [Apple]$ egrep "[bB]" myfile
2 blueflybird is me.
3 BLUEFLYBIRD IS ME.
```



grep 示例 V

④ 匹配任一字符.

- . 通常与其它元字符联合使用, 因为. 代表任一字符, 如下面的命令表示查找以任一字符开头, 第二个字符是f的行:

```
1 [Apple]$ egrep "^f" myfile
2 of course,Blueflybird is me.
3 Of Course,BLUEFLYBIRD IS ME.
```



grep 示例 VI

⑤ 匹配[]内列出的字符

- 查找包括有小写of与大写Of的行:

```
1 [Apple]$ egrep "[oO]f" myfile
2 of course,Blueflybird is me.
3 Of Course,BLUEFLYBIRD IS ME.
```

- 查找不包括小写of与大写Of的行

```
1 [Apple]$ egrep "[^oO]f" myfile
2 blueflybird is me.
3 Yes,I'm blueflybird.
4 of course,Blueflybird is me.
```



grep 示例 VII

- 查找包括符号] 和大写Y的行

```
1 [Apple]$ egrep "[Y]" myfile
2 BLUEFLYBIRD IS ME.
3 Yes,I'm blueflybird.
4 Of Course,BLUEFLYBIRD IS ME.
5 [BLUEBLUE BLUE by BLUE blue bird bird ]
```

- 查找所有以大写字母开始的行

```
1 [Apple]$ egrep "^[A-Z]" myfile
2 BLUEFLYBIRD IS ME.
3 Yes,I'm blueflybird.
4 Of Course,BLUEFLYBIRD IS ME.
```



grep 示例 VIII

⑥ \<完整单词\>

- 这里的完整单词并不是说真实的英语单词，而是指连续的字母组合中不含有数字、空格、换行符及制表符，如“myfile1”中完整单词是myfile，同样“my name1eye”中，my, name, eye各自是完整单词，而“bbeekks”无序的组合也是完整单词



grep 示例 IX

⑦ 匹配次数 `{}`、`+`、`?`、`*`

- 为了继续说明匹配次数，我们新建一个文件叫 `myfile1`，内容如下：

```
1 [Apple]$ cat myfile1
2 blue
3 blueblue
4 blueblueblue
5 BLUE
6 BLUEBLUE
7 BLUEBLUEBLUEBLUE
8 BBLLUUEE
9 BLUEEE
```



grep 示例 X

- 如果我要查找两个（精确的两个）连续BLUE相连的行

```
1 [Apple]$ egrep "(BLUE){2}" myfile1
2 BLUEBLUE
3 BLUEBLUEBLUEBLUE
```

- 注意：为了查找BLUEBLUE，用RE表示则是(BLUE){2}，而不是BLUE{2}，如果是BLUE{2}，则表示BLUEE
- 警告：本来要查BLUEBLUE，为甚麽上面还列出了三个连续的BLUE呢？因为BLUEBLUEBLUE包括了BLUEBLUE



grep 示例 XI

- 如果想查找BLUEBLUE那一行，则必须使用定位符，如：

```
1 [Apple]$ egrep "^(BLUE){2}\>" myfile1
2 BLUEBLUE
```

- 上面的例子查找以两个连续BLUE开始的完整单词的行
- 如要查找以两个或者两个以上连续BLUE开始的完整单词，则在2后面加上逗号：

```
1 [Apple]$ egrep "^(BLUE){2,}\>" myfile1
2 BLUEBLUE
3 BLUEBLUEBLUEBLUE
```



grep 示例 XII

- 如查找以一至三个连续BLUE开始的完整单词的行:

```
1 [Apple]$ egrep "(BLUE){1,3}" myfile1
2 BLUE
3 BLUEBLUE
4 BLUEBLUEBLUEBLUE
5 BLUEEEE
```

- 如果要查找以至少 1 个或多个BLUE开始的完整单词行:

```
1 [Apple]$ egrep "^(BLUE)+\>" myfile1
2 BLUE
3 BLUEBLUE
4 BLUEBLUEBLUEBLUE
```



grep 示例 XIII

- 如果要查找以 0 个或多个BLUE开始的完整单词行:

```
1 [Apple]$ egrep "^(BLUE)*\>" myfile1
2 BLUE
3 BLUEBLUE
4 BLUEBLUEBLUEBLUE
```

- 如果要查找 0 个或 1 个以BLUE开始的完整单词行:

```
1 [Apple]$ egrep "^(BLUE)?\>" myfile1
2 BLUE
```



grep 示例 XIV

- 非常值得注意的是：上面的命令 `(BLUE)?\>` 中增加了 `\>` 是否完整单词的测试，如果去掉 `\>`，那么：

```
1 [Apple]$ egrep "^(BLUE)?" myfile1
2 blue
3 blueblue
4 blueblueblue
5 BLUE
6 BLUEBLUE
7 BLUEBLUEBLUEBLUE
8 BBLLUUEE
9 BLUEEE
```



grep 示例 XV

- 由于`^(BLUE)?`表示以 0 个或 1 个`BLUE`开始的行，没有测试`BLUE`是否是完整单词，所以可以匹配任何行，而加了完整单词测试的命令，则只能匹配有 1 个`BLUE`开始并且是完整单词的行匹配。



附录：特殊符号的名字 I

符号	英文名	中文名
~	tilde	波浪号
!	exclamation	感叹号
@	at	at
#	number	井号
\$	dollar	美元符
%	percent	百分号
^	caret	control 符
&	ampersand or and	and 符
*	asterisk	星号
★	star	星号
+	plus	加号



附录：特殊符号的名字 II

-	minus	减号
-	hyphen	连字符
_	underline	下划线
—	dash	破折号
'	quote or apostrophe	单引号
"	double quote	双引号
`	back quote	反引号
.	period	英文句号
,	comma	逗号
:	colon	冒号
;	semicolon	分号
[left bracket	左方括号
]	right bracket	右方括号
(left parenthesis	左圆括号
)	right parenthesis	右圆括号



附录：特殊符号的名字 III

{	left curly bracket, or left brace	左花括号
}	right brace	右花括号
《》	French quotes	书名号
...	ellipsis	省略号
	pipe or or	管道线
	parallel	平行
/	virgule, or slash	斜线
\	backslash	反斜线
?	question mark	问号
=	equal	等号
≡	identical	恒等
≅	congruent	全等



附录：特殊符号的名字 IV



The End

The End of Chapter VI.

