

TURING 图灵程序设计丛书

[PACKT]
Publishing

[印] Sarath Lakshman 著 门佳 译

Linux Shell Scripting Cookbook

Linux Shell 脚本攻略



人民邮电出版社
POSTS & TELECOM PRESS

Linux Shell Scripting Cookbook

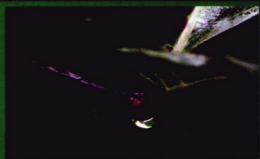
Linux Shell脚本攻略

本书通过细致剖析实际应用中的110多个案例，使许多看似复杂的Linux shell脚本任务迎刃而解

作者在Linux shell脚本方面的经验，使他能够以一种相当清晰且友好的方式有效地分享他的知识。本书会帮助读者利用少量命令的组合完成诸如文本处理、文件管理、备份等复杂的数据管理工作

本书将告诉你如何

- ◆ 利用shell命令快速开发常规任务
- ◆ 综合应用grep、find、sed和awk等常用命令
- ◆ 凭借短短几个命令行从Web挖掘数据的shell脚本
- ◆ 利用归档工具运行并自动化各种任务，诸如自动备份和存储
- ◆ 理解文件系统、文件类型以及文件管理
- ◆ 用shell创建以及维护文件或目录归档、压缩格式和加密技术
- ◆ 通过shell脚本设置以太网和无线LAN
- ◆ 使用登录技术监控网络上的各种动态



- 直截了当而便于应用的写作风格
- 常见任务和问题的精挑细选
- 针对问题而精心组织的高效解法
- 对实际操作的细致分析
- 对扩展解决方案的深入探讨

[PACKT]
PUBLISHING

图灵社区: www.it-ebooks.com.cn

反馈 投稿 推荐信箱: contact@turingbook.com

热线: (010)51095186转604



分类建议 计算机/操作系统/Linux

人民邮电出版社网址: www.ptpress.com.cn

ISBN 978-7-115-26472-5



9 787115 264725 >

ISBN 978-7-115-26472-5

定价: 49.00元

TURING 图灵程序设计丛书

[印] Sarath Lakshman 著 门佳 译

Linux Shell Scripting Cookbook

Linux Shell 脚本攻略

人民邮电出版社
北 京

图书在版编目 (CIP) 数据

Linux Shell脚本攻略 / (印) 拉克什曼
(Lakshman, S.) 著 ; 门佳译. — 北京 : 人民邮电出版
社, 2011.11

(图灵程序设计丛书)

书名原文: Linux Shell Scripting Cookbook

ISBN 978-7-115-26472-5

I. ①L… II. ①拉… ②门… III. ①
Linux操作系统—程序设计 IV. ①TP316.89

中国版本图书馆CIP数据核字(2011)第198231号

内 容 提 要

本书是 Linux Shell 编程的实战秘籍。它集合了众多适合于实战的命令行脚本攻略, 并辅以大量案例以及细致的讲解。

本书的读者不仅包括 Shell 编程的新手, 也包括那些对这一领域相当熟悉的专业人士。对新手而言, 本书的内容由浅入深且紧贴实践, 使得他们能够快速地从中学以致用, 而专业人士也能从本书中发现一些新鲜的东西, 使自己的技巧更加纯熟。

图灵程序设计丛书 Linux Shell脚本攻略

-
- ◆ 著 [印] Sarath Lakshman
 - 译 门 佳
 - 责任编辑 王军花
 - 执行编辑 王一枝
 - ◆ 人民邮电出版社出版发行 北京市崇文区夕照寺街14号
邮编 100061 电子邮件 315@ptpress.com.cn
网址 <http://www.ptpress.com.cn>
三河市潮河印业有限公司印刷
 - ◆ 开本: 800×1000 1/16
印张: 17
字数: 402千字 2011年11月第1版
印数: 1~3 000册 2011年11月河北第1次印刷
著作权合同登记号 图字: 01-2011-2969 号
ISBN 978-7-115-26472-5
-

定价: 49.00元

读者服务热线: (010)51095186转604 印装质量热线: (010)67129223

反盗版热线: (010)67171154

版 权 声 明

Copyright © 2011 Packt Publishing. First published in the English language under the title *Linux Shell Scripting Cookbook*.

Simplified Chinese-language edition copyright © 2011 by Posts & Telecom Press. All rights reserved.

本书中文简体字版由 Packt Publishing 授权人民邮电出版社独家出版。未经出版者书面许可，不得以任何方式复制或抄袭本书内容。

版权所有，侵权必究。



谨以此书献给我点亮生活方向的父母。



译者序

计算机技术类的书大致可以分为两类：入门类和实战类。入门类的书，多是从最基础、最简单的内容开始，循序渐进、由浅入深地展开，旨在引导读者完成一个从“无”到“有”的过程；而实战类的书，则是假定读者已经具备相关的背景知识，将重点放在特定情境下的实践以及技巧上，帮助读者实现从“有”到“精”的转变。得益于开源社区以及爱好者们的无私贡献，任何一位初学者都可以在Internet上轻松获取各种Linux/Unix shell相关的基础教程。而就如何运用shell来解决现实世界中的问题，相关的知识与技巧多是散落在论坛、博客之中，要么不够全面系统，要么不易查找。如果你也有此烦恼，那么不妨看看本书。从英文书名中采用的cookbook一词便可知这本书的风格如同居家菜谱（cookbook）一样直截了当：列出菜名（待解决的问题），然后给出做法（具体的命令或脚本）。市面上类似风格的书不单此一本，O'Reilly公司出版的*Unix Powertools*、*Bash Cookbook*也都采用了这种写法。从内容上看，本书并非是简单的重复，而更像是对两位“前辈”的补充，它加入了系统管理自动化、Web页面的抓取与解析、编写Gmail与twitter客户端、利用Git进行备份、结构化文本处理等较新的技巧。即便是已经看过*Unix Powertools*和*Bash Cookbook*的读者，我相信他们也能够在这本书中发现一些新东西。

本书作者是一位21岁的在校大学生，除了编写此书之外，他还是SLYNIX GNU/Linux发行版的研发人员。联想21时的自己，实在有些惭愧。

在翻译的过程中，除了修正代码或文字中出现的逻辑或印刷错误之外，对于原文的某些部分，译者并没有严格地依照字面含义翻译，而是根据原书内容，在不更改作者原意同时保证技术正确的前提下，进行了适当的调整，使其便于理解。

在豆瓣上看书评的时候，经常会看到对于某些中文版书翻译质量的抱怨，有时我也忍不住要牢骚几句。正因为如此，在翻译的过程中，我会竭力确保译文的正确性与流畅性，希望自己不会成为日后被抱怨的对象。不过和不存在没有bug的代码一样，本书中也难免有错，对此，我愿负全责。

从2011年4月接手这份翻译工作，到7月交出译稿，前后3个月的时间里少了很多陪伴家人的时间。在此要感谢我的父母和我的未婚妻姗姗，感谢你们对我的理解与支持。

前 言

GNU/Linux可谓是一款不同凡响的操作系统，它拥有一个稳定、可靠且极其强大的完备的开发环境。作为与操作系统进行沟通的原生界面（native interface），shell能够控制整个操作系统的运作。理解shell脚本可以让你更好地了解操作系统，同时还能帮助你通过短短几行脚本自动地将大部分手头工作搞定，从而节省大量的时间。shell脚本可以和许多外部命令行工具结合起来完成信息查询、简化文本处理、调度任务运行时间、生成报表以及发送邮件之类的工作。尽管不少shell命令也配有对应的文档，但是仍然不太好理解。本书收集了诸多适合于实战的命令行脚本攻略，同时辅以详细的讲解。内容上涵盖了大多数重要的Linux命令的用法，其中包括大量的实例。本书能够帮你借助几个命令来完成涉及文本处理、文件管理、备份等任务的繁杂的数据处理工作。

仅凭一行代码就能搞定复杂的文本处理任务，你想成为这样的命令行高手吗？想过写几个shell脚本和报表工具来找点儿乐子，或是做点动真格的系统管理工作吗？那么本书就是为你量身打造的。好了，开始上路吧！

本书内容

第1章涵盖了如终端打印、数学运算、数组、操作符、函数、别名、文件重定向等可以通过Bash脚本来完成的一系列初级任务。作为入门篇，本章目的在于让读者掌握Bash中的基本概念及特性。

第2章展示了GNU/Linux下多个命令在不同情境下的实战用法。介绍了cat、md5sum、find、tr、sort、uniq、split、rename、look等重要命令。本章考查了用户可能会遇到并可借鉴的各种切实可行的用例。

第3章包含了多个与文件和文件系统相关的任务攻略。本章演示了如何生成大体积文件，将文件系统写入文件并挂载，查找并删除重复文件，统计文件行数，创建ISO镜像，收集文件细节信息、符号链接操作、文件权限及属性的详情，等等。

第4章以大量实例讲解了GNU/Linux下大部分命令行文本处理工具，同时还细致描述了正则表达式及sed和awk等命令。本章在各种实例中就大多数常见的文本处理任务，详细地剖析了其解决方案。

第5章包含了多个与Internet和Web相关的shell脚本，旨在帮助读者了解如何使用shell脚本同Web打交道，从而实现采集及解析Web页面数据，以POST和GET的方式发送用户数据，编写Web

服务的客户端，下载Web页面等任务的自动化处理。

第6章结合脚本实例，演示了用于数据备份、归档、压缩等的若干命令以及用法。本章还介绍了tar、gzip、bunzip、cpio、lzma、dd、rsync、git、squashfs等命令，并讨论了一些重要的加密技术。

第7章讨论了Linux环境下的联网实践以及一些有助于编写基于网络的shell脚本的命令。为了照顾新手，本章一开头先介绍了一些网络基础知识。接下来的重头戏包括借助SSH实现无密码登录，通过网络传送文件，列出网络中的活动主机，以多播方式进行消息传播，等等。

第8章考查了Linux系统活动监视相关的实例以及日志记录和报表生成。本章讲解了诸如计算磁盘使用情况、监视用户访问、CPU占用、syslog、查看常用命令等任务。

第9章包含一系列系统管理方面的实战攻略。它介绍了用于完成系统信息采集、使用脚本进行用户管理、向用户发送消息、大图片缩放、通过shell访问MySQL数据库等任务的各种命令。

阅读本书的前提

只要具备任何一种GNU/Linux平台的一般性使用经验都有助于你更轻松地阅读本书。我们已竭尽所能地确保书中的所有例子清晰明了，并尽可能简单易懂。在Linux平台下学习的好奇心是你阅读本书所需的唯一条件。我们为你提供了循序渐进的辅导来解决书中有关脚本编写的难题。为了运行并测试书中的例子，我们推荐安装Ubuntu Linux。当然，其他的Linux发行版也足以胜任绝大多数任务。你会发现就编写shell脚本来说，本书绝对是一份通俗易懂的参考资料，同时它也是一位帮你编写出高效脚本的良师益友。

本书读者

如果你是一位初、中级用户，并且希望通过掌握快速编写脚本的技巧来完成各类事务处理，而又不愿去逐页翻阅手册，那么本书就是写给你的。你不用了解任何shell脚本或Linux的工作原理，只需要参照书中类似的例子和描述就可以动手了。对于中、高级用户以及系统管理员或程序员而言，本书则是一份绝佳的参考资料。

本书约定

本书用多种不同格式的文本来区分不同种类的信息。下面是各类格式的例子及对其所代表的含义的解说。

正文中出现的代码以如下形式显示：“我们可以通过printf来使用格式化字符串。”

代码块以如下形式显示：

```
#!/bin/bash
#Filename: printf.sh

printf "%-5s %-10s %-4s\n" No Name Mark
```

```
printf "%-5s %-10s %-4.2f\n" 1 Sarath 80.3456
printf "%-5s %-10s %-4.2f\n" 2 James 90.9989
printf "%-5s %-10s %-4.2f\n" 3 Jeff 77.564
```

任何命令行输入或输出写成如下形式：

```
$ chmod +s executable_file

# chown root.root executable_file
# chmod +s executable_file
$ ./executable_file
```



警告或重要的提示出现在这里。



建议和窍门则会以这种方式出现。

读者反馈

我们十分欢迎读者的反馈意见。我们想知道你对本书的看法：喜欢哪些部分，不喜欢哪些部分。这些反馈对于协助我们编写出真正对读者有所裨益的书至关重要。

你只需要向feedback@packtpub.com发送电子邮件，并在邮件标题中注明书名即可。

如果你确实需要并希望我们能够出版其他领域的书，请在www.packtpub.com中的SUGGEST A TITLE表格内留言，或者发送邮件至suggest@packtpub.com。

如果你在某方面有所专长并且愿意参与图书编写，请参阅我们的作者指南www.packtpub.com/authors。

客户支持

现在你已经拥有了这本由Packt出版的图书，为了让你的付出得到最大的回报，我们还为你提供了其他许多方面的服务。



下载本书的示例代码

如果你是通过<http://www.packtpub.com>的注册账户购买的图书，可以从该账户中下载你购买过的所有Packt图书中的示例代码。如果你是从其他地方购买本书的，可以访问<http://www.packtpub.com/support>并进行注册，我们将会为您发送一封附有示例代码文件的电子邮件。

勘误

尽管我们已经竭尽全力确保本书内容准确,但错误终难避免。如果你发现了书中的任何错误,无论是出现在正文还是代码中的,我们都非常乐于见到你将错误提交给我们。这样不仅能够减少其他读者的困惑,还能帮助我们改进本书后续版本的质量。如果需要提交勘误,请访问<http://www.packtpub.com/support>,选择相应的书名,单击勘误提交表格链接,就可以开始输入你的勘误信息了。一旦通过验证,我们将接受你的提交,同时勘误内容也将被上传到我们的网站,或者被添加到对应书目勘误区的现有勘误表中。任何图书当前的勘误都可以通过<http://www.packtpub.com/support>来查看。

侵权行为

各种媒体在Internet上一直饱受版权侵害的困扰。Packt坚持对版权和授权进行全力保护。如果你在Internet上发现我社图书的任何形式的盗版,请立即为我们提供地址或网站名称,以便我们采取进一步的措施。

疑难解答

如果你对本书的某方面抱有疑问,请通过questions@packtpub.com联系我们,我们会尽力为你解决。

致谢

我要感谢一直以来给予我极大支持和鼓励的好友和家人。感谢我的朋友Anu Mahadevan和Neenu Jacob,他们热情并耐心地通读了所有章节,更在此书的编写过程中提出了一些修改意见。同样要感谢Atanu Datta先生,他帮助我构思出了章节的标题。最后,非常感谢Packt出版社的工作人员,感谢他们帮助我,才使此书得以出版。



目 录

第 1 章 小试牛刀.....1	1.10 调试脚本.....22
1.1 简介.....1	1.10.1 预备知识.....22
1.2 终端打印.....3	1.10.2 实战演练.....22
1.2.1 实战演练.....3	1.10.3 补充内容.....23
1.2.2 补充内容.....4	1.11 函数和参数.....23
1.3 玩转变量和环境变量.....5	1.11.1 实战演练.....23
1.3.1 预备知识.....5	1.11.2 补充内容.....24
1.3.2 实战演练.....6	1.12 读取命令序列输出.....25
1.3.3 补充内容.....7	1.12.1 预备知识.....26
1.4 通过 shell 进行数学运算.....8	1.12.2 实战演练.....26
1.4.1 预备知识.....9	1.12.3 补充内容.....26
1.4.2 实战演练.....9	1.13 以不按回车键的方式读取字符“n”.....27
1.5 玩转文件描述符和重定向.....10	1.13.1 预备知识.....27
1.5.1 预备知识.....10	1.13.2 实战演练.....27
1.5.2 实战演练.....11	1.14 字段分隔符和迭代器.....28
1.5.3 补充内容.....13	1.14.1 预备知识.....28
1.6 数组和关联数组.....15	1.14.2 实战演练.....29
1.6.1 预备知识.....15	1.15 比较与测试.....30
1.6.2 实战演练.....15	1.15.1 预备知识.....30
1.6.3 补充内容.....16	1.15.2 实战演练.....30
1.7 使用别名.....17	第 2 章 命令之乐.....34
1.7.1 预备知识.....17	2.1 简介.....34
1.7.2 实战演练.....17	2.2 用 cat 进行拼接.....34
1.7.3 补充内容.....18	2.2.1 预备知识.....34
1.8 获取终端信息.....18	2.2.2 实战演练.....35
1.8.1 预备知识.....18	2.2.3 工作原理.....35
1.8.2 实战演练.....18	2.2.4 补充内容.....35
1.9 获取、设置日期和延时.....19	2.3 录制与回放终端会话.....37
1.9.1 预备知识.....19	2.3.1 预备知识.....37
1.9.2 实战演练.....19	2.3.2 实战演练.....37
1.9.3 补充内容.....21	

2.3.3 工作原理	37	2.13.2 工作原理	65
2.4 文件查找与文件列表	38	2.14 交互输入自动化	66
2.4.1 预备知识	38	2.14.1 预备知识	66
2.4.2 实战演练	38	2.14.2 实战演练	66
2.4.3 补充内容	39	2.14.3 工作原理	67
2.5 玩转 xargs	45	2.14.4 补充内容	67
2.5.1 预备知识	45		
2.5.2 实战演练	45	第3章 以文件之名	69
2.5.3 工作原理	46	3.1 简介	69
2.5.4 补充内容	46	3.2 生成任意大小的文件	69
2.6 用 tr 进行转换	49	3.3 文本文件的交集与差集	70
2.6.1 预备知识	49	3.3.1 预备知识	71
2.6.2 实战演练	49	3.3.2 实战演练	71
2.6.3 工作原理	49	3.4 查找并删除重复文件	73
2.6.4 补充内容	50	3.4.1 预备知识	73
2.7 校验和与核实	52	3.4.2 实战演练	73
2.7.1 预备知识	52	3.4.3 工作原理	74
2.7.2 实战演练	52	3.4.4 参考	75
2.7.3 工作原理	52	3.5 创建长路径目录	75
2.7.4 补充内容	53	3.5.1 预备知识	75
2.8 排序、单一与重复	53	3.5.2 实战演练	76
2.8.1 预备知识	54	3.6 文件权限、所有权和粘滞位	76
2.8.2 实战演练	54	3.6.1 预备知识	76
2.8.3 工作原理	54	3.6.2 实战演练	78
2.8.4 补充内容	55	3.6.3 补充内容	79
2.9 临时文件命名与随机数	58	3.7 创建不可修改文件	80
2.9.1 实战演练	58	3.7.1 预备知识	80
2.9.2 工作原理	58	3.7.2 实战演练	80
2.10 分割文件和数据	59	3.8 批量生成空白文件	80
2.10.1 工作原理	59	3.8.1 预备知识	81
2.10.2 补充内容	59	3.8.2 实战演练	81
2.11 根据扩展名切分文件名	61	3.9 查找符号链接及其指向目标	81
2.11.1 实战演练	61	3.9.1 预备知识	81
2.11.2 工作原理	61	3.9.2 实战演练	81
2.12 批量重命名和移动	63	3.10 列举文件类型统计信息	82
2.12.1 预备知识	63	3.10.1 预备知识	83
2.12.2 实战演练	63	3.10.2 实战演练	83
2.12.3 工作原理	64	3.10.3 工作原理	84
2.13 拼写检查与词典操作	65	3.11 环回文件与挂载	84
2.13.1 实战演练	65	3.11.1 预备知识	85

3.11.2 实战演练	85	4.4.1 预备知识	107
3.11.3 补充内容	86	4.4.2 实战演练	107
3.12 生成 ISO 文件及混合 ISO	87	4.4.3 补充内容	108
3.12.1 预备知识	87	4.5 统计特定文件中的词频	109
3.12.2 实战演练	87	4.5.1 预备知识	109
3.12.3 补充内容	87	4.5.2 实战演练	109
3.13 查找文件差异并进行修补	89	4.5.3 工作原理	110
3.13.1 实战演练	89	4.5.4 参考	110
3.13.2 补充内容	90	4.6 sed 入门	110
3.14 head 与 tail —— 打印文件的 前 10 行和后 10 行	90	4.6.1 实战演练	111
3.15 只列出目录的其他方法	92	4.6.2 补充内容	111
3.15.1 预备知识	93	4.7 awk 入门	113
3.15.2 实战演练	93	4.7.1 实战演练	113
3.16 在命令行中用 pushd 和 popd 快速定位	93	4.7.2 工作原理	113
3.16.1 预备知识	93	4.7.3 补充内容	114
3.16.2 实战演练	93	4.8 替换文本或文件中的字符串	117
3.16.3 补充内容	94	4.8.1 预备知识	117
3.17 统计文件的行数、单词数和字符数	95	4.8.2 实战演练	118
3.17.1 预备知识	95	4.8.3 补充内容	118
3.17.2 实战演练	95	4.8.4 参考	118
3.17.3 补充知识	95	4.9 压缩或解压缩 JavaScript	119
3.18 打印目录树	96	4.9.1 预备知识	119
3.18.1 预备知识	96	4.9.2 工作原理	119
3.18.2 实战演练	96	4.9.3 工作原理	120
3.18.3 补充内容	97	4.9.4 参考	121
第 4 章 让文本飞	98	4.10 对文件中的行、单词和字符进行 迭代	121
4.1 简介	98	4.10.1 预备知识	121
4.2 正则表达式入门	99	4.10.2 实战演练	121
4.2.1 预备知识	99	4.10.3 工作原理	122
4.2.2 实战演练	99	4.10.4 参考	122
4.2.3 工作原理	100	4.11 按列合并文件	122
4.2.4 补充内容	101	4.11.1 工作原理	122
4.3 用 grep 在文件中搜索文本	101	4.11.2 参考	123
4.3.1 预备知识	101	4.12 打印文件或行中的第 n 个单词或列	123
4.3.2 实战演练	101	4.12.1 预备知识	123
4.3.3 补充内容	103	4.12.2 实战演练	123
4.4 用 cut 按列切分文件	107	4.12.3 参考	124
		4.13 打印不同行或样式之间的文本	124
		4.13.1 预备知识	124

4.13.2 实战演练	124	第5章 一团乱麻? 没这回事	136
4.13.3 参考	125	5.1 入门	136
4.14 用脚本检验回文字符串	125	5.2 网站下载	136
4.14.1 预备知识	125	5.2.1 预备知识	136
4.14.2 工作原理	125	5.2.2 实战演练	136
4.14.3 工作原理	126	5.2.3 补充内容	137
4.14.4 补充内容	127	5.3 以格式化纯文本形式下载网页	138
4.14.5 参考	128	5.4 cURL 入门	139
4.15 以逆序形式打印行	128	5.4.1 预备知识	139
4.15.1 预备知识	128	5.4.2 实战演练	139
4.15.2 实战演练	128	5.4.3 补充内容	140
4.15.3 工作原理	129	5.4.4 参考	142
4.15.4 参考	129	5.5 从命令行访问 Gmail	142
4.16 解析文本中的电子邮件地址和 URL	129	5.5.1 实战演练	142
4.16.1 预备知识	129	5.5.2 工作原理	143
4.16.2 实战演练	129	5.5.3 参考	144
4.16.3 工作原理	130	5.6 解析网站数据	144
4.16.4 参考	130	5.6.1 实战演练	144
4.17 打印文件中某个样式之前或之后的 n 行	130	5.6.2 工作原理	144
4.17.1 预备知识	131	5.6.3 参考	145
4.17.2 实战演练	131	5.7 制作图片抓取器及下载工具	145
4.17.3 参考	132	5.7.1 实战演练	145
4.18 在文件中移除包含某个单词的句子	132	5.7.2 工作原理	146
4.18.1 预备知识	132	5.7.3 参考	147
4.18.2 实战演练	132	5.8 网页相册生成器	147
4.18.3 工作原理	133	5.8.1 预备知识	147
4.18.4 参考	133	5.8.2 实战演练	147
4.19 用 awk 实现 head、tail 和 tac	133	5.8.3 工作原理	148
4.19.1 预备知识	133	5.8.4 参考	149
4.19.2 实战演练	133	5.9 Twitter 命令行客户端	149
4.19.3 工作原理	134	5.9.1 预备知识	149
4.19.4 参考	134	5.9.2 实战演练	149
4.20 文本切片与参数操作	134	5.9.3 工作原理	150
4.20.1 实战演练	134	5.9.4 参考	150
4.20.2 参考	135	5.10 基于 Web 后端的定义查询工具	151
		5.10.1 预备知识	151
		5.10.2 实战演练	151
		5.10.3 工作原理	152
		5.10.4 参考	152

5.11 查找网站中的无效链接.....152	6.8.2 实战演练.....171
5.11.1 预备知识.....152	6.8.3 补充内容.....171
5.11.2 实战演练.....153	6.9 加密工具与散列.....172
5.11.3 工作原理.....153	6.10 用 raync 备份系统快照.....174
5.11.4 参考.....153	6.10.1 实战演练.....174
5.12 跟踪网站变更.....154	6.10.2 补充内容.....175
5.12.1 预备知识.....154	6.11 用 Git 备份版本控制.....176
5.12.2 实战演练.....154	6.11.1 预备知识.....176
5.12.3 工作原理.....155	6.11.2 实战演练.....176
5.12.4 参考.....155	6.12 用 dd 克隆磁盘.....178
5.13 以 POST 方式发送网页并读取 响应.....155	6.12.1 预备知识.....179
5.13.1 预备知识.....156	6.12.2 实战演练.....179
5.13.2 实战演练.....156	6.12.3 补充内容.....180
5.13.3 补充内容.....157	6.12.4 参考.....180
5.13.4 参考.....157	
第 6 章 B 计划.....158	第 7 章 无网不利.....181
6.1 简介.....158	7.1 简介.....181
6.2 用 tar 归档.....158	7.2 联网知识入门.....181
6.2.1 预备知识.....158	7.2.1 新手上路.....181
6.2.2 实战演练.....159	7.2.2 实战演练.....182
6.2.3 补充知识.....159	7.2.3 补充内容.....182
6.2.4 参考.....163	7.2.4 参考.....186
6.3 用 cpio 归档.....163	7.3 使用 ping.....186
6.4 用 gunzip 或 gzip 压缩.....164	7.3.1 实战演练.....186
6.4.1 实战演练.....164	7.3.2 补充内容.....187
6.4.2 补充内容.....164	7.4 列出网络上所有的活动主机.....188
6.4.3 参考.....166	7.4.1 新手上路.....188
6.5 用 bunzip 或 bzip 压缩.....166	7.4.2 实战演练.....188
6.5.1 实战演练.....166	7.4.3 工作原理.....189
6.5.2 补充内容.....167	7.4.4 补充内容.....191
6.5.3 参考.....168	7.4.5 参考.....191
6.6 用 lzma 压缩.....168	7.5 传输文件.....191
6.6.1 实战演练.....168	7.5.1 新手上路.....191
6.6.2 补充内容.....169	7.5.2 实战演练.....191
6.6.3 参考.....169	7.5.3 补充内容.....192
6.7 用 zip 归档和压缩.....169	7.5.4 参考.....194
6.8 超高压缩率的 squashfs 文件系统.....170	7.6 用脚本设置以太网与无线 LAN.....194
6.8.1 预备知识.....171	7.6.1 新手上路.....194
	7.6.2 实战演练.....194
	7.6.3 工作原理.....196

7.6.4 参考	196	8.6.3 工作原理	216
7.7 用 SSH 实现无密码自动登录	196	8.6.4 参考	217
7.8 用 SSH 在远程主机上运行命令	198	8.7 用 watch 监视命令输出	217
7.8.1 新手上路	198	8.7.1 实战演练	217
7.8.2 实战演练	198	8.7.2 补充内容	217
7.8.3 补充内容	200	8.8 对文件及目录访问进行记录	218
7.8.4 参考	200	8.8.1 新手上路	218
7.9 在本地挂载点上挂载远程驱动器	201	8.8.2 实战演练	218
7.9.1 新手上路	201	8.8.3 工作原理	218
7.9.2 实战演练	201	8.9 用 logrotate 管理日志文件	219
7.9.3 参考	201	8.9.1 新手上路	219
7.10 在网络上发送多播式窗口消息	201	8.9.2 实战演练	219
7.10.1 新手上路	201	8.10 用 syslog 记录日志	220
7.10.2 实战演练	201	8.10.1 新手上路	220
7.10.3 工作原理	202	8.10.2 实战演练	221
7.10.4 参考	203	8.10.3 参考	221
7.11 网络流量与端口分析	203	8.11 通过监视用户登录找出入侵者	221
7.11.1 新手上路	203	8.11.1 新手上路	222
7.11.2 实战演练	203	8.11.2 实战演练	222
7.11.3 补充内容	204	8.11.3 工作原理	223
第 8 章 当个好管家	205	8.12 监视远程磁盘的健康情况	224
8.1 简介	205	8.12.1 新手上路	224
8.2 统计磁盘的使用情况	205	8.12.2 实战演练	224
8.2.1 新手上路	206	8.12.3 工作原理	225
8.2.2 实战演练	206	8.12.4 参考	226
8.2.3 补充内容	206	8.13 找出系统中用户的活动时段	226
8.3 计算命令执行时间	210	8.13.1 新手上路	226
8.4 与当前登录用户、启动日志及启动故障的相关信息	212	8.13.2 实战演练	226
8.4.1 新手上路	212	8.13.3 工作原理	227
8.4.2 实战演练	212	第 9 章 管理重任	228
8.5 打印出 10 条最常使用的命令	214	9.1 简介	228
8.5.1 新手上路	214	9.2 收集进程信息	228
8.5.2 实战演练	214	9.2.1 新手上路	228
8.5.3 工作原理	215	9.2.2 实战演练	229
8.6 列出 1 小时内占用 CPU 最多的 10 个进程	215	9.2.3 补充内容	231
8.6.1 新手上路	215	9.2.4 参考	234
8.6.2 实战演练	215	9.3 杀死进程以及发送或响应信号	234
		9.3.1 新手上路	235

9.3.2 实战演练	235	9.8.3 补充内容	243
9.3.3 补充内容	235	9.9 从 Bash 中读写 MySQL 数据库	244
9.4 which、whereis、file、 what is 与平均负载	237	9.9.1 新手上路	244
9.5 向用户终端发送消息	238	9.9.2 实战演练	244
9.5.1 新手上路	239	9.9.3 工作原理	247
9.5.2 实战演练	239	9.10 用户管理脚本	248
9.5.3 工作原理	240	9.10.1 实战演练	248
9.6 收集系统信息	240	9.10.2 工作原理	249
9.7 用 /proc 收集信息	241	9.11 图像文件的批量缩放及格式转换	251
9.8 用 cron 进行调度	242	9.11.1 新手上路	251
9.8.1 新手上路	242	9.11.2 实战演练	251
9.8.2 实战演练	242	9.11.3 工作原理	253
		9.11.4 参考	254



第1章

小试牛刀

1

本章内容

- 终端打印
- 玩转变量与环境变量
- 通过shell进行数学运算
- 玩转文件描述符与重定向
- 数组和关联数组
- 使用别名
- 获取终端信息
- 获取、设置日期及延时
- 调试脚本
- 函数和参数
- 读取命令序列输出
- 以不按回车键的方式获取字符“n”
- 字段分隔符和迭代器
- 比较与测试

1.1 简介

诸多类UNIX操作系统的设计令人惊叹。即便是在数十年之后的今天，UNIX风格的操作系统架构仍是有史以来的最佳设计之一。这种架构最重要的一个特性就是命令行界面或shell。shell环境使得用户能与操作系统的内核进行交互操作。术语“脚本”更多涉及的是这类环境。编写脚本通常使用某种基于解释器的编程语言。而shell脚本不过就是一些文件，我们能将一系列需要执行的命令写入其中，然后通过shell来执行这些脚本。

本书使用Bash（Bourne Again Shell），它是目前大多数GNU/Linux系统默认的shell环境。由于GNU/Linux是基于UNIX风格架构的最杰出的操作系统，书中大部分案例和讨论都假定是在Linux系统环境下进行的。

本章的主要目的是让读者了解shell环境并熟悉shell的基本特性。命令都是在shell终端中输入并执行。打开终端后，就会出现一个提示符。其形式通常如下：

```
username@hostname$
```

或者

```
root@hostname#
```

要么就简单地以\$或#表示。

\$表示普通用户，#表示超级用户（root user）。超级用户是Linux系统中权限最高的用户。

shell脚本通常是一个以`#!/`起始的文本文件，如下所示：

```
#!/bin/bash
```

Linux环境下的任何脚本语言，都是以这样一个被称为shebang^①的特殊行作为起始的。在这行中，字符`#!/`被置于解释器路径之前。`/bin/bash`是Bash的路径。

有两种运行脚本的方式。一种是将脚本作为`sh`的命令行参数，另一种是将脚本作为具有执行权限的可执行文件。

将脚本作为命令行参数时的运行方式如下：

```
$ sh script.sh # 假设脚本位于当前目录下
```

或者

```
$ sh /home/path/script.sh # 使用script.sh的完整路径
```

如果将脚本作为`sh`的命令行参数来运行，那么脚本中的shebang行也就没什么用处了。

为了使shell脚本能够自己独立运行，需要具备可执行权限。要使脚本独立运行，必须利用shebang行。它通过使用位于`#!/`之后的解释器来运行脚本。至于脚本的可执行权限，可以通过以下方式设置：

```
$ chmod a+x script.sh
```

该命令赋予所有用户`script.sh`文件的可执行权限。这个脚本能以下列方式执行：

```
$ ./script.sh #./ 表示当前目录
```

或者

```
$ /home/path/script.sh # 使用脚本的完整路径
```

shell程序读取脚本的首行，查看shebang行是否为`#!/bin/bash`。它会识别`/bin/bash`，并在内部以如下命令行执行该脚本：

```
$ /bin/bash script.sh
```

当打开一个终端的时候，该终端最初会执行一组命令来定义诸如提示文本、颜色等各类设置。这组命令来自于用户`home`目录中的`.bashrc`脚本文件（`~/.bashrc`）。Bash还维护了一个历史记录文件`~/.bash_history`，用于保存用户运行过的命令。`~`是一种简写，代表用户`home`目录的路径。

在Bash中，每个命令或是命令序列是通过使用分号或换行符来分隔的。比如：

```
$ cmd1 ; cmd2
```

它等同于：

```
$ cmd1
$ cmd2
```

① shebang这个词其实是两个字符名称的组合。在Unix的行话里，用sharp或hash（有时候是mesh）来称呼字符“#”，用bang来称呼惊叹号“!”，因而shebang合起来就代表了这两个字符。详情请参考：[http://en.wikipedia.org/wiki/Shebang_\(Unix\)](http://en.wikipedia.org/wiki/Shebang_(Unix))。——译者注（书中所有的注均为译者注。）

字符#指明注释的开始。注释部分以#为起始，一直延续到行尾。注释行通常用于为代码提供注释信息，或者用于暂停执行某行代码。^①

现在让我们继续。

1.2 终端打印

终端作为交互式工具，用户可以通过它与shell环境进行交互。在终端中打印文本是绝大多数shell脚本和工具日常需要进行的基本任务。能够执行打印的方法有很多，格式也各有不同。

1.2.1 实战演练

echo是用于终端打印的基本命令。

在默认情况下，echo在每次调用后会添加一个换行符。

```
$ echo "Welcome to Bash"
Welcome to Bash
```

只需要使用带双引号的文本，结合echo命令就可以将该文本在终端中打印出来。类似地，不带双引号的文本也可以得到同样的输出结果：

```
$ echo Welcome to Bash
Welcome to Bash
```

使用单引号也可以完成同样的任务：

```
$ echo 'text in quote'
```

这些方法看起来相似，但各有一些特殊用途和副作用。思考下面这行命令：

```
$ echo "cannot include exclamation - ! within double quotes"
```

这条命令将会返回：

```
bash: !: event not found error
```

因此，如果你希望打印!，那就不要将其放入双引号中，或者你可以在其之前加上一个特殊的转义字符(\)将!转义。

```
$ echo Hello world !
```

或者

```
$ echo 'Hello world !'
```

或者

```
$ echo "Hello world \!" #Escape character \ prefixed.
```

当在echo中使用带双引号的文本时，你应该在echo之前使用set +H，以便能够正常地显示!。每种方法的副作用如下：

^① shell不执行脚本代码中的任何注释部分。

- 使用不带引号的echo时,你没法在所要显示的文本中使用,因为在bash shell中被用作命令定界符。
- 以echo hello;hello为例,echo hello被视为一个命令,第二个hello则被视为另一个命令。
- 使用带单引号的echo时,Bash不会对单引号中的变量(如\$var)求值,而只是照原样显示。

这意味着: \$ echo '\$var'将会返回 \$var, 而\$ echo \$var将会根据变量\$var定义与否,返回\$var的值,或者什么都不返回。

另一个可用于终端打印的命令是printf。printf使用的参数和C语言中的printf函数一样。例如:

```
$printf "Hello world"
```

printf使用引用文本或由空格分隔的参数。我们可以在printf中使用格式化字符串。我们还可以指定字符串的宽度、左右对齐方式等。在默认情况下,printf并不像echo命令一样会自动添加换行符,我们必须在需要的时候手动添加,比如在下面的脚本中:

```
#!/bin/bash
#文件名: printf.sh

printf "%-5s %-10s %-4s\n" No Name Mark
printf "%-5s %-10s %-4.2f\n" 1 Sarath 80.3456
printf "%-5s %-10s %-4.2f\n" 2 James 90.9989
printf "%-5s %-10s %-4.2f\n" 3 Jeff 77.564
```

我们会得到如下格式化的输出:

No	Name	Mark
1	Sarath	80.35
2	James	91.00
3	Jeff	77.56

%s、%c、%d和%f都是格式替代符(format substitution character),其所对应的参数可以置于带引号的格式字符串之后。

%-5s指明了一个格式为左对齐且宽度为5的字符串替代(-表示左对齐)。如果不用-指定对齐方式,字符串则采用右对齐形式。宽度指定了保留给某个变量的字符数。对Name而言,保留宽度是10。因此,任何Name字段的内容都会被显示在10字符宽的保留区域内,如果内容不足10字符,余下的则以空格符填充。

对于浮点数,我们可以使用其他参数对小数部分进行舍入。

对于Mark字段,我们将其格式化为%-4.2f,其中.2指定保留2个小数位。注意,在每行格式字符串后都有一个换行符\n。

1.2.2 补充内容

一定要留神的是echo和printf中的标志(如-e、-n等)应该出现在命令行内任何字符串之前,

否则Bash会将其视为另外一个字符串。

1. 在echo中转义换行符

在默认情况下，echo会将一个换行符追加到输出文本的尾部。可以使用标志-n来忽略结尾的换行符。echo同样接受双引号字符串内的转义序列（escape sequence）作为参数。如果需要使用转义序列，则采用echo-e'包含转义序列的字符串'这种形式。例如：

```
echo -e "1\t2\t3"
123
```

2. 打印彩色输出

在终端中生成彩色输出相当好玩，我们可以使用转义序列来实现。

每种颜色都有对应的颜色码。比如：重置=0，黑色=30，红色=31，绿色=32，黄色=33，蓝色=34，洋红=35，青色=36，白色=37。

要打印彩色文本，可输入如下命令：

```
echo -e "\e[1;31m This is red text \e[0m"
```

\e[1;31m将颜色设为红色，\e[0m将颜色重新置回。你只需要将31替换成想要的颜色码就可以了。

要设置彩色背景，经常使用的颜色码是：重置=0，黑色=40，红色=41，绿色=42，黄色=43，蓝色=44，洋红=45，青色=46，白色=47。

要打印彩色文本，可输入如下命令：

```
echo -e "\e[1;42m Green Background \e[0m"
```

1.3 玩转变量和环境变量

变量是任何一种编程语言必不可少的组成部分，用于存放各类数据。脚本语言通常不需要在使用变量之前声明其类型。只需要直接赋值就可以了。在Bash中，每一个变量的值都是字符串。无论你给变量赋值时有没有使用引号，值都会以字符串的形式存储。有一些特殊的变量会被shell环境和操作系统环境用来存储一些特别的值，这类变量就被称为环境变量。

让我们来看一些实例。

1.3.1 预备知识

变量采用常见的命名方式进行命名。当一个应用程序执行的时候，它接收一组环境变量。可以使用env命令在终端中查看所有与此终端进程相关的环境变量。对于每个进程，在其运行时的环境变量可以使用下面的命令来查看：

```
cat /proc/$PID/environ
```

其中，将PID设置成相关进程的进程ID（PID总是一个整数）。

假设有一个叫做gedit的应用程序正在运行。我们可以使用pgrep命令获得gedit的进程ID：

环境变量是未在当前进程中定义，而从父进程中继承而来的变量。例如环境变量HTTP_PROXY，它定义了一个Internet连接应该使用哪一个代理服务器。

该环境变量通常被设置成：

```
HTTP_PROXY=http://192.168.0.2:3128
export HTTP_PROXY
```

export命令用来设置环境变量。至此之后，从当前shell脚本执行的任何程序都会继承这个变量。我们可以按照自己的需要，在执行的应用程序或者shell脚本中导出特定的变量。在默认情况下，有很多标准环境变量可供shell使用。

PATH就是其中之一。通常，变量PATH包含：

```
$ echo $PATH
```

```
/home/slynnux/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/
sbin:/bin:/usr/games
```

在给出所要执行的命令后，shell自动在PATH环境变量所包含的目录列表中（各目录路径之间以冒号分隔）查找对应的可执行文件。\$PATH通常定义在/etc/environment或/etc/profile或~/.bashrc中。如果需要PATH中添加一条新路径，可以使用：

```
export PATH="$PATH:/home/user/bin"
```

也可以使用

```
$ PATH="$PATH:/home/user/bin"
$ export PATH
```

```
$ echo $PATH
/home/slynnux/bin:/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/
sbin:/bin:/usr/games:/home/user/bin
```

这样，我们就将/home/user/bin添加到了PATH中。

还有一些众所周知的环境变量：HOME、PWD、USER、UID、SHELL等。

1.3.3 补充内容

让我们再多看些有关普通变量和环境变量的技巧。

1. 获得字符串长度

可以用下面的方法获得变量值的长度：

```
length=${#var}
```

例如：

```
$ var=12345678901234567890
$ echo ${#var}
20
```

length就是字符串所包含的字符数。



2. 识别当前的shell版本

可以用下面的方法获知当前使用的是哪种shell:

```
echo $SHELL
```

也可以用

```
echo $0
```

例如:

```
$ echo $SHELL
/bin/bash
```

```
$ echo $0
bash
```

3. 检查是否为超级用户

UID是一个重要的环境变量, 可以用于检查当前脚本是以超级用户还是以普通用户的身份运行的。例如:

```
if [ $UID -ne 0 ]; then
echo Non root user. Please run as root.
else
echo "Root user"
fi
```

root用户的UID是0。

4. 修改Bash提示字符串 (username@hostname:~\$)

当我们打开一个终端或是运行一个shell, 都会看到类似于user@hostname: /home/\$的提示字符串。不同GNU/Linux发布版中的提示及颜色也略有不同。我们可以利用PS1环境变量来定制提示文本。默认的shell提示文本是在文件~/.bashrc中的某一行设置的。

□ 可以使用如下命令列出设置PS1的那一行:

```
$ cat ~/.bashrc | grep PS1
PS1='${debian_chroot:+($debian_chroot)}\u@\h:\w$ '
```

□ 如果要设置提示字符串, 可以输入:

```
slynux@localhost: ~$ PS1="PROMPT>"
PROMPT> Type commands here # 提示字符串已经改变
```

□ 我们可以利用类似\@[1;31的特定转义序列来设置彩色的提示字符串(参考1.2节的内容)。还有一些特殊的字符可以扩展成系统参数。例如:\u可以扩展为用户名,\h可以扩展为主机名, 而\w可以扩展为当前工作目录。

1.4 通过 shell 进行数学运算

无论哪种编程语言都少不了算数操作, Bash shell同样提供了多种此类操作。

1.4.1 预备知识

在 Bash shell 环境中，可以利用 `let`、`(())` 和 `[]` 执行基本的算数操作。而在进行高级操作时，`expr` 和 `bc` 这两个工具也会非常有用。

1.4.2 实战演练

可以用普通的变量赋值方法定义数值，这时，它会被存储为字符串。然而，我们可以用一些方法使它像数字一样进行处理。

```
#!/bin/bash
no1=4;
no2=5;
```

`let` 命令可以直接执行基本的算数操作。

当使用 `let` 时，变量名之前不需要再添加 `$`，例如：

```
let result=no1+no2
echo $result
```

□ 自加操作

```
$ let no1++
```

□ 自减操作

```
$ let no1--
```

□ 简写形式

```
let no+=6
let no-=6
```

它们分别等同于 `let no=no+6` 和 `let no=no-6`。

□ 其他方法

操作符 `[]` 的使用方法和 `let` 命令类似：

```
result=${ no1 + no2 }
```

在 `[]` 中也可以使用 `$` 前缀，例如：

```
result=${ $no1 + 5 }
```

也可以使用 `(())`，但使用 `(())` 时，变量名之前需要加上 `$`：

```
result=$(( no1 + 50 ))
```

`expr` 同样可以用于基本算数操作：

```
result=`expr 3 + 4`
result=$(expr $no1 + 5)
```

以上这些方法只能用于整数运算，而不支持浮点数。

`bc` 是一个用于数学运算的高级工具，这个精密计算器包含了大量的选项。我们可以借助

它执行浮点数运算并应用一些高级函数：

```
echo "4 * 0.56" | bc
2.24

no=54;
result=`echo "$no * 1.5" | bc`
echo $result
81.0
```

其他参数可以置于要执行的具体操作之前，同时以分号作为定界符，通过stdin传递给bc。

- **设定小数精度(数值范围)：**在下面的例子中，参数scale=2将小数位数设置为2。因此，bc将会输出包含两个小数位的数值。

```
echo "scale=2;3/8" | bc
0.37
```

- **进制转换：**用bc可以将一种进制系统转换为另一种。来看看如何将十进制转换成二进制，然后再将二进制转换回十进制：

```
#!/bin/bash
#用途：数字转换

no=100
echo "obase=2;$no" | bc
1100100
no=1100100
echo "obase=10;ibase=2;$no" | bc
100
```

- **计算平方以及平方根：**

```
echo "sqrt(100)" | bc #Square root
echo "10^10" | bc #Square
```

1.5 玩转文件描述符和重定向

文件描述符是与文件输入、输出相关联的整数。它们用来跟踪已打开的文件。最常见的文件描述符是stdin、stdout和stderr。我们可以将某个文件描述符的内容重定向到另一个文件描述符中。下面给出一些对文件描述符进行操作和重定向的例子。

1.5.1 预备知识

我们在编写脚本的时候会频繁使用标准输入(stdin)、标准输出(stdout)和标准错误(stderr)。通过内容过滤将输出重定向到文件是我们从事的基础任务之一。当命令输出文本的时候，这些输出文本有可能是错误信息，也可能是正常的（非错误的）输出信息。单靠查看输出的文本本身，我们没法区分哪些是正常的输出文本，哪些是错误文本。不过，我们可以通过文件描述符来解决这个问题，将那些与特定描述符关联的文本提取出来。

文件描述符是与一个打开的文件或数据流相关联的整数。文件描述符0、1以及2是系统预留的。

- 0 —— stdin (标准输入)。
- 1 —— stdout (标准输出)。
- 2 —— stderr (标准错误)。

1.5.2 实战演练

用下面的方法可以将输出文本重定向或保存到一个文件中：

```
$ echo "This is a sample text 1" > temp.txt
```

这种方法通过截取文件的方式，将输出文本存储到文件temp.txt中，也就是说在把echo命令的输出写入文件之前，temp.txt中的内容首先会被清空。

接下来，再看另一个例子：

```
$ echo "This is sample text 2" >> temp.txt
```

这种方法会将文本追加到目标文件中。

>和>>并不相同。尽管这两个操作符都可以将文本重定向到文件，但是前者会先清空文件，再写入内容；而后者会将内容追加到现有文件的尾部。

可以用下面的方法查看文件内容：

```
$ cat temp.txt
This is sample text 1
This is sample text 2
```

当使用重定向操作符时，重定向的内容不会出现在终端，而是直接被导入文件。重定向操作符默认使用标准输出。如果想使用特定的文件描述符，你必须将描述符置于操作符之前。

>等同于1>；对于>>来说，情况也类似（即>>等同于1>>）。

来看看什么是标准错误以及如何对它重定向。当命令输出错误信息时，stderr信息就会被打印出来。考虑下面的例子：

```
$ ls +
ls: cannot access +: No such file or directory
```

这里，+ 是一个非法参数，因此将返回错误信息。

成功和不成功的命令



当一个命令发生错误并返回时，它会返回一个非0的退出状态；而当命令成功完成后，它会返回数字0。退出状态可以从特殊变量\$? 中获得（在命令执行语句之后立刻运行echo \$?，就可以打印出退出状态）。

下面的命令会将stderr文本打印到屏幕上，而不是文件中。

```
$ ls + > out.txt
ls: cannot access +: No such file or directory
```

然而在下方的命令中，`stdout`没有任何输出，因此会生出空文件`out.txt`。

```
$ ls + 2> out.txt # 正常运行
```

你可以将`stderr`单独重定向到一个文件，将`stdout`重定向到另一个文件：

```
$ cmd 2>stderr.txt 1>stdout.txt
```

还可以利用下面的方法将`stderr`转换成`stdout`，使得`stderr`和`stdout`都被重定向到同一个文件中：

```
$ cmd 2>&1 output.txt
```

或者采用下列方法：

```
$ cmd &> output.txt
```

有时候，在输出中可能包含一些不必要的信息（比如除错信息）。如果你不想让终端中充斥着有关`stderr`的繁枝末节，那么你可以将`stderr`的输出重定向到 `/dev/null`，保证一切都会被清除得干干净净。假设我们有三个文件，分别是`a1`、`a2`、`a3`。但是普通用户对文件`a1`没有“读-写-执行”权限。如果你需要打印文件名以`a`起始的所有文件的内容，你可以使用`cat`命令。

设置一些测试文件：

```
$ echo a1 > a1
$ cp a1 a2 ; cp a2 a3;
$ chmod 000 a1 #清除所有权限
```

尽管可以使用通配符（`*`）显示所有的文件内容，但是系统会显示一个出错信息，因为对文件`a1`没有可读权限。

```
$ cat a*
cat: a1: Permission denied
a1
a1
```

其中，`cat: a1: Permission denied`属于`stderr`。我们可以将`stderr`信息重定向到一个文件中，而`stdout`仍然保持不变。考虑如下代码：

```
$ cat a* 2> err.txt #stderr被重定向到err.txt
a1
a1
```

```
$ cat err.txt
cat: a1: Permission denied
```

观察下面的代码：

```
$ some_command 2> /dev/null
```

在这个示例中，来自`stderr`的输出被丢到文件`/dev/null`中。`/dev/null`是一个特殊的设备文件，这个文件接收到的任何数据都会被丢弃。因此，`null`设备通常也被称为位桶（bit bucket）或黑洞。

当对`stderr`或`stdout`进行重定向时，重定向的文本将传入文件。因为文本已经被重定向到文件中，也就没剩下什么东西可以通过管道（`|`）传给接下来的命令，而这些命令是通过`stdin`

来接收文本的。

但是有一个巧妙的方法可以一方面将数据重定向到文件，另一方面还可以提供一份重定向数据的副本作为后续命令的stdin。这一切都可以使用tee来实现。举个例子：要在终端中打印stdout，同时将它重定向到一个文件中，那么可以这样使用tee：

```
command | tee FILE1 FILE2
```

在下面的代码中，tee命令接收到来自stdin的数据。它将stdout的一份副本写入文件out.txt，同时将另一份副本作为后续命令的stdin。命令cat -n将从stdin中接收到的每一行数据前加上行号并写入stdout：

```
$ cat a* | tee out.txt | cat -n
cat: a1: Permission denied
1a1
2a1
```

查看out.txt的内容：

```
$ cat out.txt
a1
a1
```

注意，cat: a1: Permission denied并没有在文件内容中出现。这是因为这些信息属于stderr，而tee只能从stdin中进行读取。

默认情况下，tee命令会将文件覆盖，但它提供了一个-a选项，可以用于追加内容。例如：

```
$ cat a* | tee -a out.txt | cat -n.
```

带有参数的命令可以写成：command FILE1 FILE2依次类推，或者简简单单地用command FILE。

我们可以使用stdin作为命令参数。只需要将-作为命令的文件名参数即可：

```
$ cmd1 | cmd2 | cmd -
```

例如：

```
$ echo who is this | tee -
who is this
who is this
```

或者我们也可以将/dev/stdin作为输出文件名来使用stdin。

类似地，使用/dev/stderr代表标准错误，dev/stdout代表标准输出。这些特殊的设备文件分别对应stdin、stderr和stdout。

1.5.3 补充内容

从stdin读取输入的命令能以多种方式接收数据。另外，还可以用cat和管道来制定我们自己的文件描述符，例如：

```
$ cat file | cmd
$ cmd1 | cmd2
```

1. 将文件重定向到命令

借助重定向，我们可以像使用stdin那样从文件中读取数据：

```
$ cmd < file
```

2. 重定向脚本内部的文本块

有时候，我们需要对文本块（多行文本）像标准输入一样进行重定向。考虑一个特殊情况：源文本就位于shell脚本中。一个实用的例子是向log文件中写入头部数据，可以按照下面的方法完成：

```
#!/bin/bash
cat <<EOF>log.txt
LOG FILE HEADER
This is a test log file
Function: System statistics
EOF
```

在cat <<EOF>log.txt与下一个EOF行之间的所有文本行都会被当做stdin数据。log.txt文件的内容打印如下：

```
$ cat log.txt
LOG FILE HEADER
This is a test log file
Function: System statistics
```

3. 自定义文件描述符

文件描述符是用于访问文件的一个抽象指针。存取文件离不开被称为“文件描述符”的特殊数字。0、1和2分别是stdin、stdout和stderr的预留描述符。

我们可以使用exec命令创建自定义的文件描述符。如果你对用其他编程语言进行文件编程非常熟悉，你可能已经注意到了文件打开模式。通常来说，会使用3种模式。

- ☐ 只读模式。
- ☐ 截断模式。
- ☐ 追加模式。

< 操作符用于从文件中读取至stdin。> 操作符用于截断模式的文件写入（数据在目标文件内容被截断之后写入）。>> 操作符用于追加模式的文件写入（数据被添加到文件的现有内容中，而且该目标文件中原有的内容不会丢失）。文件描述符可以用以上三种模式中的任意一种来创建。

为读取文件创建一个文件描述符：

```
$ exec 3<input.txt # 使用文件描述符3打开并读取文件
```

我们可以这样使用它：

```
$ echo this is a test line > input.txt
$ exec 3<input.txt
```

现在你就可以在命令中使用文件描述符3了。例如：

```
$ cat <&3
this is a test line
```

如果要再次读取，我们就不能再继续使用文件描述符3了，而是需要用exec重新分配文件描述符3以便用于读取。

创建一个文件描述符用于写入（截断模式）：

```
$ exec 4>output.txt # 打开文件用于写入
```

例如：

```
$ exec 4>output.txt
$ echo newline >&4
$ cat output.txt
newline
```

创建一个文件描述符用于写入（追加模式）：

```
$ exec 5>>input.txt
```

例如：

```
$ exec 5>>input.txt
$ echo appended line >&5
$ cat input.txt
newline
appended line
```

1.6 数组和关联数组

数组是shell脚本非常重要的组成部分，它借助索引将多个独立的数据存储为一个集合。

1.6.1 预备知识

Bash同时支持普通数组和关联数组。普通数组只能使用整数作为数组索引，而关联数组可以使用字符串作为数组索引。

关联数组在很多操作中相当有用。Bash从4.0版本开始支持关联数组。也就是说，使用旧版本的Bash则享受不到这种便利。

1.6.2 实战演练

定义数组的方法有很多种。可以在单行中使用一系列值来定义一个数组：

```
array_var=(1 2 3 4 5 6)
#这些值将会存储在以0为起始索引的连续位置上
```

另外，还可以将数组定义成一组索引-值（index-value pair）：

```
array_var[0]="test1"
array_var[1]="test2"
array_var[2]="test3"
array_var[3]="test4"
array_var[4]="test5"
array_var[5]="test6"
```

打印出特定索引的数组元素内容：

```
$ echo ${array_var[0]}
test1

index=5
$ echo ${array_var[$index]}

test6
```

以清单形式打印出数组中的所有值：

```
$ echo ${array_var[*]}
test1 test2 test3 test4 test5 test6
```

你也可以使用：

```
$ echo ${array_var[@]}
test1 test2 test3 test4 test5 test6
```

打印数组长度（即数组中元素的个数）：

```
$ echo ${#array_var[*]}
6
```

1.6.3 补充内容

关联数组从Bash 4.0版本开始被引入。借助散列技术，它成为解决很多问题的有力工具。接下来就让我们一探究竟。

1. 定义关联数组

在关联数组中，我们可以用任意的文本作为数组索引。而在普通数组中，只能用整数作为数组索引。

首先，需要使用单独的声明语句将一个变量名声明为关联数组。声明语句如下：

```
$ declare -A ass_array
```

声明之后，可以用两种方法将元素添加到关联数组中。

(1) 利用内嵌索引-值列表法，提供一个索引-值列表：

```
$ ass_array=( [index1]=val1 [index2]=val2 )
```

(2) 使用独立的索引-值进行赋值：

```
$ ass_array[index1]=val1
$ ass_array[index2]=val2
```

举个例子，试想如何用关联数组为水果制订价格：

```
$ declare -A fruits_value
$ fruits_value=([apple]='100dollars' [orange]='150 dollars')
```

用下面的方法显示数组内容：

```
$ echo "Apple costs ${fruits_value[apple]}"
Apple costs 100 dollars
```



2. 列出数组索引

每一个数组元素都有一个索引用于查找。普通数组和关联数组具有不同的索引类型。我们可以用下面的方法获取数组的索引列表：

```
$ echo ${!array_var[*]}
```

也可以使用：

```
$ echo ${!array_var[@]}
```

以先前提到的`fruits_value`数组为例，运行如下：

```
$ echo ${!fruits_value[*]}
orange apple
```

对于普通数组，这个方法同样可行。

1.7 使用别名

别名就是一种便捷方式，以省去用户输入一长串命令序列的麻烦。

1.7.1 预备知识

别名有多种实现方式，可以使用函数，也可以使用`alias`命令。

1.7.2 实战演练

可以按照下面的方式创建一个别名：

```
$ alias new_command='command sequence'
```

为安装命令`apt-get install`创建别名：

```
$ alias install='sudo apt-get install'
```

这样一来，我们就可以用`install pidgin`代替`sudo apt-get install pidgin`了。

`alias`命令的作用只是暂时的。一旦关闭当前终端，所有设置过的别名就失效了。为了使别名设置一直保持作用，可以将它放入`~/.bashrc`文件中。因为每当一个新的shell进程生成时，都会执行`~/.bashrc`中的命令。

```
$ echo 'alias cmd="command seq"' >> ~/.bashrc
```

如果需要删除别名，只用将其对应的语句从`~/.bashrc`中删除，或者使用`unalias`命令。

另一种创建别名的方法是定义一个具有新名称的函数，并把它写入`~/.bashrc`。

我们可以创建一个别名`rm`，它能够删除原始文件，同时在`backup`目录中保留副本：

```
alias rm='cp $@ ~/.backup; rm $@'
```

当你创建别名时，如果已经有同名的别名存在，那么原有的别名设置将被新的取代。

1.7.3 补充内容

有时别名也会造成安全问题。下面来看看应该如何识别这些隐患。

对别名进行转义

`alias`命令能够为任何重要的命令创建别名，不过你可能未必总是希望使用这些别名。我们可以将所要运行的命令进行转义，从而忽略当前定义过的所有别名。例如：

```
$ \command
```

字符\`\`对命令实施转义，使我们可以执行原本的命令，而不是这些命令的别名替身。在不信任的环境下执行特权命令，通过在命令前加上\`\`来忽略可能存在的别名设置总是一个不错的安全实践。因为攻击者可能已经利用别名将某些特权命令替换成了一些别有用心命令，借此来盗取用户输入的重要信息。

1.8 获取终端信息

编写命令行shell脚本的时候，总是免不了大量处理当前终端的相关信息，比如行数、列数、光标位置和遮盖密码字段等。这则攻略将帮助你学习如何收集和处理终端设置。

1.8.1 预备知识

`tput`和`stty`是两款终端处理工具。下面来看看如何用它们完成各种不同的任务。

1.8.2 实战演练

获取终端的行数和列数：

```
tput cols  
tput lines
```

打印出当前终端名：

```
tput longname
```

将光标移动到方位(100,100)处：

```
tput cup 100 100
```

设置终端背景色：

```
tput setb no
```

其中，`no`可以在0到7之间取值。

将文本前景色设置为白色：

```
tput setf no
```

其中，`no`可以在0到7之间取值。

设置文本样式为粗体：



```
tput bold
```

设置下划线的起止:

```
tput smul
```

```
tput rmul
```

删除当前光标位置到行尾的所有内容:

```
tput ed
```

在输入密码的时候, 不能让输入的内容显示出来。在下面的例子中, 我们将看到如何使用 `stty` 来实现这一要求:

```
#!/bin/sh
#Filename: password.sh
echo -e "Enter password: "
stty -echo
read password
stty echo
echo
echo Password read.
```

其中, 选项 `-echo` 禁止将输出发送到终端, 而选项 `echo` 则允许发送输出。

1.9 获取、设置日期和延时

很多应用程序需要以不同的格式打印日期, 设置日期和时间, 以及根据日期和时间执行操作。延时通常用于在程序执行过程中提供一段等待时间(比如1秒钟)。比如在每隔五秒钟执行一次监视任务的这类脚本中, 则需要掌握如何在程序中加入延时。这则攻略会告诉你怎么处理日期以及延时。

1.9.1 预备知识

我们能够以多种格式打印日期, 也可以在命令行中设置日期。在类UNIX系统中, 日期被存储为一个整数, 其大小为自世界标准时间^①1970年1月1日0时0分0秒^②起所流逝的秒数。这种计时方式称之为纪元时或UNIX时间。来看看如何读取和设置它们吧。

1.9.2 实战演练

读取日期:

```
$ date
Thu May 20 23:09:04 IST 2010
```

打印纪元时:

-
- ① UTC (Coordinated Universal Time), 又称世界标准时间或世界协调时间。UTC是以原子时秒长为基础, 在时刻上尽量接近于世界时的一种时间计量系统。
 - ② UNIX认为UTC 1970年1月1日0点是纪元时间。POSIX标准推出后, 这个时间也被称为POSIX时间。

```
$ date +%s
1290047248
```

纪元时被定义为从世界标准时间1970年1月1日0时0分0秒起至当前时刻的总秒数，不包括闰秒^①。当计算两个日期或两段时间的差值时，纪元时很有用处。你可以很容易得出两个特定时间戳的纪元时，并计算出两者之间的差值，由此就能知道两个日期之间相隔了多少秒。

我们可以从给定格式的日期串中得出对应的纪年时。你在输入时有多种日期格式可供选择。如果你要从系统日志中或者其他任何标准应用程序生成的输出中获取日期，那就完全不用烦心日期格式的问题。要将日期串转换成纪元时，只需要输入如下代码：

```
$ date --date "Thu Nov 18 08:07:21 IST 2010" +%s
1290047841
```

选项 `--date` 用于提供日期串作为输入。但我们可以使用任意的日期格式化选项来打印输出。将日期串作为输入能够用来获知给定的日期是星期几。

例如：

```
$ date --date "Jan 20 2001" +%A
Saturday
```

表1-1是一份日期格式字符串列表。

表 1-1

日期内容	格 式
星期	%a (例如: Sat)
	%A (例如: Saturday)
月	%b (例如: Nov)
	%B (例如: November)
日	%d (例如: 31)
固定格式日期 (mm/dd/yy)	%D (例如: 10/18/10)
年	%y (例如: 10)
	%Y (例如: 2010)
小时	%I 或 %H (例如: 08)
分钟	%M (例如: 33)
秒	%S (例如: 10)
纳秒	%N (例如: 695208515)
UNIX纪元时 (以秒为单位)	%s (例如: 1290049486)

^① 国际原子时的误差为每日数纳秒，而世界时的误差为每日数毫秒。针对这种情况，一种称为世界标准时间的折衷时标于1972年面世。为确保UTC与世界时相差不会超过0.9秒，在有需要情况下会在UTC内加上正或负闰秒，因此UTC与国际原子时之间会出现若干整数秒的差别。位于巴黎的国际地球自转事务中央局负责决定何时加入闰秒。

用格式串结合 + 作为date命令的参数，可以按照你的选择打印出对应格式的日期。例如：

```
$ date "+%d %B %Y"
20 May 2010
```

设置日期和时间：

```
# date -s "格式化的日期字符串"
```

例如：

```
# date -s "21 June 2009 11:01:22"
```

有时候，我们需要检查一组命令所花费的时间，那就可以采用下面的方式：

```
#!/bin/bash
#文件名: time_take.sh
start=$(date +%s)
commands;
statements;

end=$(date +%s)
difference=$(( end - start))
echo Time taken to execute commands is $difference seconds.
```

另一种方法则是使用timescriptpath来得到执行脚本所使用的时间。

1.9.3 补充内容

编写以循环方式运行的监视脚本时，设置时间间隔是必不可少的。让我们来看看如何生成延时。

在脚本中生成延时

为了在脚本中推迟执行一段时间，可以使用sleep：

```
$ sleep no_of_seconds.
```

例如，下面的脚本就使用tput和sleep从0开始计数到40：

```
#!/bin/bash
#Filename: sleep.sh
echo -n Count:
tput sc

count=0;
while true;
do
if [ $count -lt 40 ];
then let count++;
sleep 1;
tput rc
tput ed
echo -n $count;
else exit 0;
fi
done
```



在上面的例子中，变量count初始化为0，随后每循环一次便增加1。echo语句打印出count的值。我们用tput sc存储光标位置。在每次循环中，我们通过恢复之前存储的光标位置，在终端中打印出新的count值。恢复光标位置的命令是tput rc。tput ed清除从当前光标位置到行尾之间的所有内容，使得旧的count值可以被清除并写入新值。循环内的1秒钟延时是通过sleep命令来实现的。

1.10 调试脚本

调试功能是每一种编程语言都应该实现的重要特性之一，当出现一些始料未及的情况时，用它来生成脚本运行信息。调试信息可以帮你弄清楚是什么原因使得程序发生崩溃或行为异常。每位系统程序员都应该了解Bash提供的那些调试选项，另外还需要了解一些调试技巧。

1.10.1 预备知识

调试shell脚本不需要什么特殊工具。Bash本身就包含了一些选项，能够打印出脚本接受的参数和输入。让我们来看看该怎么做。

1.10.2 实战演练

使用选项-x，启动跟踪调试shell脚本：

```
$ bash -x script.sh
```

运行带有-x标志的脚本能打印出所执行的每一行命令以及当前状态。注意，你也可以使用sh -x script。

-x标识将脚本中执行过的每一行都输出到stdout。不过，我们也可以要求只关注脚本某些部分的命令及参数的打印输出。针对这种情况，可以在脚本中使用set built-in来启用或禁止调试打印。

- ❑ set -x: 在执行时显示参数和命令。
- ❑ set +x: 禁止调试。
- ❑ set -v: 当命令进行读取时显示输入。
- ❑ set +v: 禁止打印输入。

例如：

```
#!/bin/bash
#文件名: debug.sh
for i in {1..6}
do
    set -x
    echo $i
    set +x
done
echo "Script executed"
```



在上面的脚本中，仅在 `-x` 和 `+x` 所限制的区域内，`echo $i` 的调试信息才会被打印出来。

这种调试方法是由 Bash 的内建功能提供的。它们通常以固定的格式生成调试信息。但是在很多情况下，我们需要以自定义格式显示调试信息。这可以通过传递 `_DEBUG` 环境变量来建立这类调试风格。

请看下面的代码：

```
#!/bin/bash
function DEBUG()
{
  [ "$_DEBUG" == "on" ] && $@ || :
}

for i in {1..10}
do
  DEBUG echo $i
done
```

可以将调试功能置为 `"on"` 来运行上面的脚本：

```
$ _DEBUG=on ./script.sh
```

我们在每一个需要打印调试信息的语句前加上 `DEBUG`。如果没有把 `_DEBUG=on` 传递给脚本，那么调试信息就不会打印出来。在 Bash 中，命令 `:` 告诉 shell 不要进行任何操作。

1.10.3 补充内容

还有其他脚本调试的便捷方法，我们甚至可以巧妙地利用 shebang 来进行调试。

shebang 的妙用

把 shebang 从 `#!/bin/bash` 改成 `#!/bin/bash -xv`，这样一来，不用任何其他选项就可以启用调试功能了。

1.11 函数和参数

和其他脚本语言一样，Bash 同样支持函数。让我们看看它是如何定义和使用函数的。

1.11.1 实战演练

定义函数：

```
function fname()
{
  statements;
}
```

或者

```
fname()
{
  statements;
}
```

只需要使用函数名就可以调用某个函数：

```
$ fname ; # 执行函数
```

参数可以传递给函数，并由脚本进行访问：

```
fname arg1 arg2 ; # 传递参数
```

以下是函数fname的定义。在函数fname中，包含了各种访问函数参数的方法。

```
fname()
{
    echo $1, $2; # 访问参数1和参数2
    echo "$@";   # 以列表的方式一次性打印所有参数
    echo "$*";   # 类似于$@，但是参数被作为单个实体
    return 0;    # 返回值
}
```

类似地，参数可以传递给脚本并通过script:\$0（脚本名）访问。

- ☐ \$1是第一个参数。
- ☐ \$2是第二个参数。
- ☐ \$n是第n个参数。
- ☐ "\$@" 被扩展成 "\$1" "\$2" "\$3"等。
- ☐ "\$*" 被扩展成 "\$1c\$2c\$3"，其中c是IFS的第一个字符。
- ☐ "\$@" 用得最多。由于 "\$*" 将所有的参数当做单个字符串，因此它很少被使用。

1.11.2 补充内容

让我们再研究Bash函数的一些技巧。

1. 递归函数

在Bash中，函数同样支持递归（可以调用自身的函数）。例如，F() { echo \$1; F hello; sleep 1; }。

Fork炸弹

```
:(){ !:;& }::
```



这个递归函数能够调用自身，不断地生成新的进程，最终造成拒绝服务攻击。函数调用前的 & 将子进程放入后台。这段危险的代码会分支出大量的进程，因而被称为Fork炸弹。

上面这段代码要理解起来可不容易。请参阅维基百科http://en.wikipedia.org/wiki/Fork_bomb，那里列出了有关Fork炸弹的细节以及更详细的解释。

可以通过修改配置文件/etc/security/limits.conf来限制可生成的最大进程数来避开这枚炸弹。

2. 导出函数

函数也能像环境变量一样用`export`导出，如此一来，函数的作用域就可以扩展到子进程中，如下：

```
export -f fname
```

3. 读取命令返回值（状态）

我们可以按照下面的方式获取命令或函数的返回值：

```
cmd;  
echo $?;
```

`$?` 会给出命令`cmd`的返回值。

返回值被称为退出状态。它可用于分析命令执行成功与否。如果命令成功退出，那么退出状态为0，否则为非0。

我们可以按照下面的方法检测某个命令是否成功结束：

```
#!/bin/bash  
#文件名: success_test.sh  
CMD="command" # command指代你要检测退出与否的目标命令  
status  
$CMD  
if [ $? -eq 0 ];  
then  
echo "$CMD executed successfully"  
else  
echo "$CMD terminated unsuccessfully"  
fi
```

4. 向命令传递参数

命令的参数能够以不同的格式进行传递。假设`-p`、`-v`是可用选项，`-k NO`是另一个可以接受数字的选项，同时该命令还接受一个文件名作为参数，那么，它有如下几种执行方式：

```
$ command -p -v -k 1 file
```

或者

```
$ command -pv -k 1 file
```

或者

```
$ command -vpk 1 file
```

或者

```
$ command file -pvk 1
```

1.12 读取命令序列输出

shell脚本最棒的特性之一就是可以轻松地将多个命令或工具组合起来生成输出。一个命令的输出可以作为另一个命令的输入，而这个命令的输出又会传递至另一个命令，依次类推。这种

命令组合的输出可以被存储在一个变量中。这则攻略将演示如何组合多个命令以及如何读取其输出。

1.12.1 预备知识

输入通常是通过`stdin`或参数传递给命令。输出要么出现在`stderr`，要么出现在`stdout`。当我们组合多个命令时，同时将`stdin`用于输入，`stdout`用于输出。

这些命令被称为过滤器（filter）。我们使用管道（pipe）来连接每一个过滤器。管道操作符是“|”。例如：

```
$ cmd1 | cmd2 | cmd3
```

这里我们组合了三个命令。`cmd1`的输出传递给`cmd2`，而`cmd2`的输出传递给`cmd3`，最终（来自`cmd3`）的输出将会被打印或导入某个文件。

1.12.2 实战演练

请看下面的代码：

```
$ ls | cat -n > out.txt
```

`ls`的输出（当前目录内容的列表）被传给`cat -n`，`cat -n`为通过`stdin`所接收到输入内容加上行号，然后将输出重定向到文件`out.txt`。

我们可以用下面的方法读取命令序列的输出：

```
cmd_output=$(COMMANDS)
```

这种方法也被称为子shell（subshell）。例如：

```
cmd_output=$(ls | cat -n)
echo $cmd_output
```

另一种被称为反引用（back-quote）的方法也可以用于存储命令输出：

```
cmd_output=`COMMANDS`
```

例如：

```
cmd_output=`ls | cat -n`
echo $cmd_output
```

反引用与单引号可不是一回事，它位于键盘的~键上。

1.12.3 补充内容

有很多种方法可以给命令分组。来看看其中的几种。

1. 利用子shell生成一个独立的进程

子shell本身就是独立的进程。可以使用`()`操作符来定义一个子shell；

```
pwd;  
(cd /bin; ls);  
pwd;
```

当命令在子shell中执行时，不会对当前shell有任何影响，所有的改变仅限于子shell内。例如，当用cd命令改变子shell的当前目录时，这种变化不会反映到主shell环境中。

pwd命令打印出工作目录的路径。

cd命令将当前目录更改为给定的目录路径。

2. 通过引用子shell的方式保留空格和换行符

假设我们使用子shell或反引用的方法将命令的输出读入一个变量中，可以将其放入双引号中，以保留空格和换行符（\n）。例如：

```
$ cat text.txt  
1  
2  
3  
  
$ out=$(cat text.txt)  
$ echo $out  
1 2 3 # 丢失了换行符 \n  
$ out="$(cat text.txt)"  
$ echo $out  
1  
2  
3
```

1.13 以不按回车键的方式读取字符“n”

read是一个重要的Bash命令，用于从键盘或标准输入中读取文本。我们可以使用read以交互的形式读取来自用户的输入，不过read能做的可远不止这些。接着我们来演示read命令的一些重要选项。

1.13.1 预备知识

任何一种编程语言的输入库大多都是从键盘读取输入，但只有当回车键按下时，才标志着输入完毕。在有些重要情形下是没法按回车键的，输入结束与否是基于字符数或某个特定字符来决定的。例如，在一个游戏中，当按下+键时，小球就会向上移动。那么若每次都要按下+键，然后再按回车键来确认已经按下+键，这就显然太低效了。read命令提供了一种不需要按回车键就能够搞定这个任务的方法。

1.13.2 实战演练

下面的语句从输入中读取n个字符并存入变量variable_name；

```
read -n number_of_chars variable_name
```

例如：

```
$ read -n 2 var
$ echo $var
```

read还有很多其他选项。让我们来看看它们。

用不回显 (non-echoed) 的方式读取密码：

```
read -s var
```

显示提示信息：

```
read -p "Enter input:" var
```

在特定时限内读取输入：

```
read -t timeout var
```

例如：

```
$ read -t 2 var
#在2秒钟内将键入的字符串读入变量var
```

用定界符结束输入行：

```
read -d delim_charvar
```

例如：

```
$ read -d ":" var
hello: #var被设置为hello
```

1.14 字段分隔符和迭代器

内部字段分隔符 (Internal Field Separator, IFS) 是shell脚本中的一个重要概念。在处理文本数据时，它可是相当有用。我们将会讨论把单个数据流划分成不同数据元素的定界符。内部字段分隔符是用于特定用途的定界符。IFS是存储定界符的环境变量。它是当前shell环境使用的默认定界字符串。

考虑一种情形：我们需要迭代一个字符串或CSV (Comma Separated Value, 逗号分隔型数值) 中的单词。在前者中，我们使用IFS="."；在后者中，则使用IFS=","。让我们看看应该怎么做。

1.14.1 预备知识

考虑CSV数据的情况：

```
data="name,sex,rollno,location"
#我们可以使用IFS读取变量中的每一个条目
oldIFS=$IFS
IFS=,now,
for item in $data;
```

```
do
echo Item: $item
done
```

```
IFS=$oldIFS
```

输入如下:

```
Item: name
Item: sex
Item: rollno
Item: location
```

IFS的默认值为空白字符(换行符、制表符或者空格)。

当IFS被设置为逗号时, shell将逗号解释成一个定界符, 因此变量 \$item在每次迭代中读取由逗号分隔的字符串作为变量值。

如果没有把IFS设置成",", 那么上面的脚本会将全部数据作为单个字符串打印出来。

1.14.2 实战演练

我们以 /etc/passwd 为例, 看看IFS的另一种用法。在文件 /etc/passwd 中, 每一行包含了由冒号划分的多个条目。文件中的每行都对应一位用户的相关属性。

考虑这样的输入: root:x:0:0:root:/root:/bin/bash。每行的最后一项指定了用户的默认shell。可以按照下面的方法巧妙地利用IFS打印出用户以及他们默认的shell:

```
#!/bin/bash
#用途: 演示IFS的用法
line="root:x:0:0:root:/root:/bin/bash"
oldIFS=$IFS;
IFS=":"
count=0
for item in $line;
do

[ $count -eq 0 ] && user=$item;
[ $count -eq 6 ] && shell=$item;
let count++;
done;
IFS=$oldIFS
echo $user's shell is $shell;
```

输出为:

```
root's shell is /bin/bash
```

对一系列值进行迭代的时候, 循环非常有用。Bash提供了多种类型的循环。下面来看看怎样使用它们。

for循环

```
for var in list;
do
```

```

commands; # 使用变量$var
done
list can be a string, or a sequence.

```

我们可以轻松地生成不同的序列。

echo {1..50}能够生成一个从1到50的数字列表。

echo {a..z}或{A..Z}，或是使用{a..h}生成部分列表。类似地，将这些方法结合起来，我们就可以连接（concatenate）数据。

下面的代码中，变量i在每次迭代的过程里都会保存一个字符，范围从a到z：

```

for i in {a..z}; do actions; done;

for循环也可以采用C语言中for循环的格式。例如：

for ((i=0;i<10;i++))
{
  commands; # 使用变量$i
}

```

while循环

```

while condition
do
  commands;
done

```

用true作为循环条件能够产生无限循环。

until循环

在Bash中还可以使用一个特殊的循环until。它会一直执行循环直到给定的条件为真。例如：

```

x=0;
until [ $x -eq 9 ]; # [ $x -eq 9 ] is the condition
do let x++; echo $x;
done

```

1.15 比较与测试

程序中的流程控制是由比较和测试语句来处理的。Bash同样具备多种与UNIX系统级特性相兼容的执行测试的方法。

1.15.1 预备知识

我们可以用if、if else以及逻辑运算符来执行测试，而用一些比较运算符来比较数据项。另外，有一个test命令也可以用来进行测试。让我们来看看如何使用这些命令。

1.15.2 实战演练

if条件：

```

if condition;
then
  commands;
fi

else if和else:

if condition;
then
  commands;
elif condition;
then
  commands
else
  commands
fi

```

if和else语句可以进行嵌套。if的条件判断部分可能会变得很长，但可以用逻辑运算符将它变得简洁一些：

```

[ condition ] && action; # 如果condition为真，则执行action
[ condition ] || action; # 如果condition为假，则执行action

```

&&是逻辑与运算符，||是逻辑或运算符。写Bash脚本时，这是一个很有用的技巧。现在来了解一下条件和比较操作。

算术比较

条件通常被放置在封闭的中括号内。一定要注意在[或]与操作数之间有一个空格。如果忘记了这个空格，脚本就会报错。例如：

```
[ $var -eq 0 ] or [ $var -eq 0 ]
```

对变量或值进行算术条件判断：

```

[ $var -eq 0 ] # 当 $var 等于 0 时，返回真
[ $var -ne 0 ] # 当 $var 为非 0 时，返回真

```

其他重要的操作符如下所示。

- -gt：大于。
- -lt：小于。
- -ge：大于或等于。
- -le：小于或等于。

可以按照下面的方法结合多个条件进行测试：

```

[ $var1 -ne 0 -a $var2 -gt 2 ] # 使用逻辑与-a
[ $var -ne 0 -o var2 -gt 2 ] # 逻辑或 -o

```

文件系统相关测试

我们可以使用不同的条件标志测试不同的文件系统相关属性。

- [-f \$file_var]：如果给定的变量包含正常的文件路径或文件名，则返回真。
- [-x \$var]：如果给定的变量包含的文件可执行，则返回真。

- [-d \$var]: 如果给定的变量包含的是目录, 则返回真。
- [-e \$var]: 如果给定的变量包含的文件存在, 则返回真。
- [-c \$var]: 如果给定的变量包含的是一个字符设备文件的路径, 则返回真。
- [-b \$var]: 如果给定的变量包含的是一个块设备文件的路径, 则返回真。
- [-w \$var]: 如果给定的变量包含的文件可写, 则返回真。
- [-r \$var]: 如果给定的变量包含的文件可读, 则返回真。
- [-L \$var]: 如果给定的变量包含的是一个符号链接, 则返回真。

使用方法举例如下:

```

fpath="/etc/passwd"
if [ -e $fpath ]; then
echo File exists;
else
echo Does not exist;
fi

```

字符串比较

使用字符串比较时, 最好用双中括号, 因为有时候采用单个中括号会产生错误, 所以最好避开它们。

可以检查两个字符串, 看看它们是否相同。

- [[\$str1 = \$str2]]; 当str1等于str2时, 返回真。也就是说, str1和str2包含的文本是一模一样的。

- [[\$str1 == \$str2]]; 这是检查字符串是否相等的另一种写法。

也可以检查两个字符串是否不同。

- [[\$str1 != \$str2]]; 如果str1和str2不相同, 则返回真。

我们还可以检查字符串的字母序情况, 具体如下所示。

- [[\$str1 > \$str2]]; 如果str1的字母序比str2大, 则返回真。
- [[\$str1 < \$str2]]; 如果str1的字母序比str2小, 则返回真。
- [[-z \$str1]]; 如果str1包含的是空字符串, 则返回真。
- [[-n \$str1]]; 如果str1包含的是非空字符串, 则返回真。



注意在=前后各有一个空格。如果忘记加空格, 那就不是比较关系了, 而变成了赋值语句。

使用逻辑运算符 && 和 || 能够很容易地将多个条件组合起来:

```

if [[ -n $str1 ]] && [[ -z $str2 ]];
then
commands;
fi

```

例如:

```
str1="Not empty "  
str2=""  
if [[ -n $str1 ]] && [[ -z $str2 ]];  
then  
echo str1 is non-empty and str2 is empty string.  
fi
```

输出如下:

str1 is non-empty and str2 is empty string.

test命令可以用来执行条件检测。用test有助于避免使用过多的括号。之前讲过的[]中的测试条件同样可以用于test命令。

例如:

```
if [ $var -eq 0 ]; then echo "True"; fi  
can be written as  
if test $var -eq 0 ; then echo "True"; fi
```



本章内容

- 用cat进行拼接
- 录制与回放终端会话
- 文件查找与文件列表
- 将命令输出作为命令参数 (xargs)
- 用tr进行转换
- 校验和与核实
- 排序、单一与重复
- 临时文件命名与随机数
- 分割文件和数据
- 根据扩展名切分文件名
- 用rename和mv批量重命名文件
- 拼写检查与词典操作
- 交互输入自动化

2.1 简介

各种命令可谓类UNIX系统中优美的部分。它们能帮助我们搞定各种繁杂的任务，使我们的工作变得更轻松。当你将这些命令付诸实践的时候，你肯定会爱上它们。在很多情形下它们都会让你情不自禁地大呼“wow!”一旦你尝试过Linux提供的这些利器，你一定会感到惊讶：以前没有这些命令的时候，自己是怎么熬过来的。在这些使我生活变得更舒适、工作更给力的命令中，我最钟爱的是grep、awk、sed和find。

UNIX/Linux命令行的使用是一门艺术。实践得越多，收益就越大。本章将为你介绍一些最有趣同时也是最实用的命令。

2.2 用 cat 进行拼接

cat是命令行玩家首先必须学习的命令之一。cat命令简约却不失优美。它通常用于读取、显示或拼接文件内容，不过cat所具备的能力远不止这些。

2.2.1 预备知识

用一个单行命令组合标准输入的数据与文件数据，这可是个让人挠头的难题。通常的解决方

法是将stdin重定向到一个文件，然后再将两个文件拼接到一起。但是我们可以使用cat命令一次性完成任务。

2.2.2 实战演练

cat命令是一个日常经常会使用到的简单命令。cat本身表示concatenate（拼接）。

用cat读取文件内容的一般写法是：

```
$ cat file1 file2 file3 ...
```

这个命令将作为命令行参数的文件内容拼接在一起作为输出。例如：

```
$ cat file.txt
This is a line inside file.txt
This is the second line inside file.txt
```

2.2.3 工作原理

cat包含了大量功能。让我们来看一些cat的用法。

cat命令不仅可以读取文件并拼接数据，它还能够从标准输入中进行读取。

要从标准输入中读取，就要使用管道操作符：

```
OUTPUT_FROM_SOME_COMMANDS | cat
```

类似地，我们可以用cat将输入文件的内容与标准输入拼接在一起。方法如下：

```
$ echo 'Text through stdin' | cat - file.txt
```

在上面的代码中，- 被作为来自stdin文本的文件名。

2.2.4 补充内容

cat命令另外还有一些选项可用于查看文件。来看看它们的用法。

1. 压缩空白行

出于可读性或是别的一些原因，有时文本中的多个空行需要被压缩成单个。可以用下面的方法压缩文本文件中连续的空白行：

```
$ cat -s file
```

例如：

```
$ cat multi_blanks.txt
line 1
```

```
line2
```

```
line3
```

```
line4
$ cat -s multi_blanks.txt # 压缩连续的空白符
line 1
```

```
line2
```

```
line3
```

```
line4
```

我们也可以用`tr`移除空白行:

```
$ cat multi_blanks.txt | tr -s '\n'
line 1
line2
line3
line4
```

在`tr`的这种用法中, 它将连续多个`'\n'`字符压缩成单个`'\n'` (换行符)。

2. 将制表符显示为 ^I

单从视觉上很难将制表符同连续的空格区分开。而在用Python之类的语言编写程序时, 将制表符和空格用于代码缩进, 具有特殊含义, 并进行区别对待。因此, 若应该在应该使用空格的地方误用了制表符的话, 就会产生缩进错误。仅仅在文本编辑器中进行查看是很难发现这种错误的。`cat`有一个特性, 可以将制表符重点标记出来。该特性对排除缩进错误非常有用。用`cat`命令的 `-T` 选项能够将制表符标记成 `^I`。例如:

```
$ cat file.py
def function():
    var = 5
        next = 6
    third = 7

$ cat -T file.py
def function():
^Ivar = 5
        next = 6
^Ithird = 7^I
```

3. 行号

使用`cat`命令的 `-n` 选项会在输出的每一行内容之前加上行号。别担心, `cat`命令绝不会修改你的文件, 它只是根据用户提供的选项在`stdout`生成一个修改过的输出而已。例如:

```
$ cat lines.txt
line
line
line

$ cat -n lines.txt
1 line
2 line
3 line
```

2.3 录制与回放终端会话

当你需要为别人在终端上演示某些操作或是需要准备一个命令行教程时,通常你要一边手动输入命令一边演示,或者,你也可以录制一段屏幕演示视频,然后再回放出来。如果我们将输入命令后发生的一切按照先后次序记录下来,再进行回放,从而使得观众好像身临其境一般,这个想法听起来如何?命令的输出会显示在终端上,一直到回放内容播放完毕。听起来好玩吧?所有这些都可以用script和scriptreplay命令来实现。

2.3.1 预备知识

script和scriptreplay命令在绝大多数GNU/Linux发行版上都可以找到。把终端会话记录到一个文件中会很有意思。你可以通过录制终端会话来制作命令行技巧视频教程,也可以与他人分享会话记录文件,共同研究如何使用这些命令行完成某项任务。

2.3.2 实战演练

开始录制终端会话:

```
$ script -t 2> timing.log -a output.session
type commands;
...
..
exit
```

两个配置文件被当做script命令的参数。其中一个文件(timing.log)用于存储时序信息,描述每一个命令在何时运行;另一个文件(output.session)用于存储命令输出。-t选项用于将时序数据导入stderr。2>则用于将stderr重定向到timing.log。

借助这两个文件:timing.log(存储时序信息)和output.session(存储命令输入信息),我们可以按照下面的方法回放命令执行过程:

```
$ scriptreplay timing.log output.session # 按播放命令序列输出
```

2.3.3 工作原理

通常,我们会录制桌面环境视频来作为教程使用。不过要注意的是,视频需要大量的存储空间,而终端脚本文件仅仅是一个文本文件,其文件大小不过是KB级别。

script命令同样可以用于建立可在多个用户之间进行广播的视频会话。这是件很有意思的事。来看看它是如何实现的吧。

打开两个终端,Terminal1和Terminal2。

(1)在Terminal1中输入以下命令:

```
$ mkfifo scriptfifo
```

(2) 在Terminal2中输入以下命令:

```
$ cat scriptfifo
```

(3) 返回Terminal1, 输入以下命令:

```
$ script -f scriptfifo  
$ commands;
```

如果需要结束会话, 输入`exit`并按回车键。你会得到如下信息: "Script done, file is scriptfifo"。现在, Terminal1就成为了广播员, 而Terminal2则成为了听众。

不管你在Terminal1中输入什么内容, 它都会在Terminal2或者使用了下列命令的任何终端中实时播放:

```
cat scriptfifo
```

当需要为计算机实验室或Internet上的用户群演示教程的话, 不妨考虑这个方法。它在节省带宽的同时也提供了实时体验。

2.4 文件查找与文件列表

`find`是UNIX/Linux命令行工具箱中最棒的工具之一。这个命令对编写shell脚本很有帮助, 但是多数人由于对它缺乏认识, 并不能有效地使用它。这则攻略讨论了`find`的大多数用法以及如何用它解决各种难度的问题。

2.4.1 预备知识

`find`命令的工作方式如下: 沿着文件层次结构向下遍历, 匹配符合条件的文件, 并执行相应的操作。下面来看看`find`命令的各种用例和基本用法。

2.4.2 实战演练

要列出当前目录及子目录下所有的文件和文件夹, 可以采用下面的写法:

```
$ find base_path
```

`base_path`可以是任何位置 (例如 `/home/slynux`), `find`会从该位置开始向下查找。

例如:

```
$ find . -print
```

打印文件和目录的列表

· 指定当前目录, .. 指定父目录。这是UNIX文件系统中的约定用法。

`-print`指明打印出匹配文件的文件名 (路径)。当使用 `-print` 时, '\n' 作为用于分隔文件的定界符。

`-print0`指明使用 '\0' 作为定界符来打印每一个匹配的文件名。当文件名中包含换行符时, 这个方法就有用武之地了。

2.4.3 补充内容

至此，我们已经学习了find命令最常见的用法。作为一个强大的命令行工具，find命令包含了诸多值得留意的选项。接下来让我们来看find命令的一些其他选项。

1. 根据文件名或正则表达式匹配搜索

选项-name的参数指定了文件名所必须匹配的字符串。我们可以将通配符作为参数使用。`*.txt`能够匹配所有以.txt结尾的文件名。选项-print在终端中打印出符合条件（例如-name）的文件名或文件路径，这些匹配条件作为find命令的选项给出。

```
$ find /home/slynux -name "*.txt" -print
```

find命令有一个选项-iname（忽略字母大小写），该选项的作用和-name类似，只不过在匹配名字的时候会忽略大小写。

例如：

```
$ ls
example.txt EXAMPLE.txt file.txt
$ find . -iname "example*" -print
./example.txt
./EXAMPLE.txt
```

如果想匹配多个条件中的一个，可以采用OR条件操作：

```
$ ls
new.txt some.jpg text.pdf
$ find . \( -name "*.txt" -o -name "*.pdf" \) -print
./text.pdf
./new.txt
```

上面的代码会打印出所有的.txt和.pdf文件，是因为这个find命令能够匹配所有这两类文件。`\(以及\)`用于将-name "*.txt" -o -name "*.pdf"视为一个整体。

选项-path的参数可以使用通配符来匹配文件路径或文件。-name总是用给定的文件名进行匹配。-path则将文件路径作为一个整体进行匹配。例如：

```
$ find /home/users -path "*/slynux*" -print
This will match files as following paths.
/home/users/list/slynux.txt
/home/users/slynux/eg.css
```

选项-regex的参数和-path的类似，只不过-regex是基于正则表达式来匹配文件路径的。

正则表达式是通配符匹配的高级形式，它可以指定文本模式。我们借助这种模式来匹配并打印文本。使用正则表达式进行文本匹配的一个典型例子就是从一堆文本中解析出所有的E-mail地址。一个E-mail地址通常采用name@host.root这种形式，所以，可以将其一般化为a-z0-9]+@[a-z0-9]+.[a-z0-9]+。+指明在它之前的字符类中的字符可以出现一次或多次。

下面的命令匹配.py或.sh文件：

```
$ ls
new.PY next.jpg test.py
$ find . -regex ".*\\(\\.py\\|\\.sh\\)$"
./test.py
```

类似地，`-iregex`用于忽略正则表达式的大小写。例如：

```
$ find . -iregex ".*\\(\\.py\\|\\.sh\\)$"
./test.py
./new.PY
```

2. 否定参数

`find`也可以用“`!`”否定参数的含义。例如：

```
$ find . ! -name "*.txt" -print
```

上面的`find`命令能够匹配所有不以`.txt`结尾的文件名。下面就是这个命令的运行结果：

```
$ ls
list.txt new.PY new.txt next.jpg test.py
$ find . ! -name "*.txt" -print
.
./next.jpg
./test.py
./new.PY
```

3. 基于目录深度的搜索

`find`命令在使用时会遍历所有的子目录。我们可以采用一些深度参数来限制`find`命令遍历的深度。`-maxdepth`和`-mindepth`就是这类参数。

大多数情况下，我们只需要在当前目录中进行搜索，无须再继续向下查找。对于这种情况，我们使用深度参数来限制`find`命令向下查找的深度。如果只允许`find`在当前目录中查找，深度可以设置为1；当需要向下两级时，深度可以设置为2；其他情况可以依次类推。

我们可以用`-maxdepth`参数指定最大深度。与此相似，我们也可以指定一个最小的深度，告诉`find`应该从此处开始向下查找。如果我们想从第二级目录开始搜索，那么使用`-mindepth`参数设置最小深度。使用下列命令将`find`命令向下的最大深度限制为1：

```
$ find . -maxdepth 1 -type f -print
```

该命令只列出当前目录下的所有普通文件。即使有子目录，也不会被打印或遍历。与之类似，`-maxdepth 2`最多向下遍历两级子目录。

`-mindepth`类似于`-maxdepth`，不过它设置的是`find`遍历的最小深度。这个选项可以用来查找并打印那些距离起始路径超过一定深度的所有文件。例如，打印出深度距离当前目录至少两个子目录的所有文件：

```
$ find . -mindepth 2 -type f -print
./dir1/dir2/file1
./dir3/dir4/f2
```

即使当前目录或`dir1`和`dir3`中包含有文件，它们也不会被打印出来。



`-maxdepth`和`-mindepth`应该作为`find`的第3个参数出现。如果作为第4个或之后的参数,就可能会影响到`find`的效率,因为它不得不进行一些不必要的检查。例如,如果`-maxdepth`作为第4个参数,`-type`作为第三个参数,`find`首先会找出符合`-type`的所有文件,然后在所有匹配的文件中再找出符合指定深度的那些。但是如果反过来,目录深度作为第三个参数,`-type`作为第四个参数,那么`find`就能够在找到所有符合指定深度的文件后,再检查这些文件的类型,这才是最有效的搜索顺序。

2

4. 根据文件类型搜索

类UNIX系统将一切都视为文件。文件具有不同的类型,例如普通文件、目录、字符设备、块设备、符号链接、硬链接、套接字以及FIFO等。

`-type`可以对文件搜索进行过滤。借助这个选项,我们可以为`find`命令指明特定的文件匹配类型。

只列出所有的目录:

```
$ find . -type d -print
```

将文件和目录分别列出可不是个容易事。不过有了`find`就好办了。例如,只列出普通文件:

```
$ find . -type f -print
```

只列出符号链接:

```
$ find . -type l -print
```

你可以按照表2-1列出的内容用`type`参数来匹配所需要的文件类型。

表 2-1

文件类型	类型参数
普通文件	<code>f</code>
符号链接	<code>l</code>
目录	<code>d</code>
字符设备	<code>c</code>
块设备	<code>b</code>
套接字	<code>s</code>
Fifo	<code>p</code>

5. 根据文件时间进行搜索

UNIX/Linux文件系统中的每一个文件都有三种时间戳 (timestamp), 如下所示。

- ☐ 访问时间 (`-atime`): 用户最近一次访问文件的时间。
- ☐ 修改时间 (`-mtime`): 文件内容最后一次被修改的时间。
- ☐ 变化时间 (`-ctime`): 文件元数据 (metadata, 例如权限或所有权) 最后一次改变的时间。

在UNIX中并没有所谓“创建时间”的概念。

-atime、-mtime、-ctime可作为find的时间参数。它们可以整数值给出，单位是天。这些整数值通常还带有-或+：-表示小于，+表示大于，如下所示。

□ 打印出在最近七天内被访问过的所有文件：

```
$ find . -type f -atime -7 -print
```

□ 打印出恰好在七天前被访问过的所有文件：

```
$ find . -type f -atime 7 -print
```

□ 打印出访问时间超过七天的所有文件：

```
$ find . -type f -atime +7 -print
```

类似地，我们可以根据修改时间，用-mtime进行搜索，也可以根据变化时间，用-ctime进行搜索。

-atime、-mtime以及-ctime都是基于时间的参数，其计量单位是“天”。还有其他一些基于时间的参数是以分钟作为计量单位的。这些参数包括：

□ -amin（访问时间），

□ -mmin（修改时间），

□ -cmin（变化时间）。

举例如下。

打印出访问时间超过7分钟的所有文件：

```
$ find . -type f -amin +7 -print
```

find另一个漂亮的特性是-newer参数。使用-newer，我们可以指定一个用于比较时间戳的参考文件，然后找出比参考文件更新的（更长的修改时间）所有文件。

例如，找出比file.txt修改时间更长的所有文件：

```
$ find . -type f -newer file.txt -print
```

find命令的时间戳操作处理选项对编写系统备份和维护脚本很有帮助。

6. 基于文件大小的搜索

根据文件的大小，可以这样搜索：

```
$ find . -type f -size +2k
# 大于2KB的文件
```

```
$ find . -type f -size -2k
# 小于2KB的文件
```

```
$ find . -type f -size 2k
# 大小等于2KB的文件
```

除了k之外，还可以用其他文件大小单元。

□ b——块（512字节）。

□ c——字节。



- w —— 字 (2字节)。
- k —— 千字节。
- M —— 兆字节。
- G —— 吉字节。

7. 删除匹配的文件

-delete可以用来删除find查找到的匹配文件。

删除当前目录下所有的.swp文件:

```
$ find . -type f -name "*.swp" -delete
```

8. 基于文件权限和所有权的匹配

文件匹配可以根据文件权限进行。列出具有特定权限的所有文件:

```
$ find . -type f -perm 644 -print
# 打印出权限为644的文件
```

以Apache Web服务器为例。Web服务器上的PHP文件需要具有合适的执行权限。我们可以用下面的方法找出那些没有设置好执行权限的PHP文件:

```
$ find . -type f -name "*.php" ! -perm 644 -print
```

也可以根据文件的所有权进行搜索。用选项 -user USER就能够找出由某个特定用户所拥有的文件:

参数USER可以是用户名或UID。

例如, 打印出用户slynux拥有的所有文件:

```
$ find . -type f -user slynux -print
```

9. 结合find执行命令或动作

find命令可以借助选项 -exec与其他命名进行结合。-exec算得上是find最强大的特性之一。看看应如何使用 -exec选项。

在前一节中, 我们用 -perm找出了所有权限不当的PHP文件。这次的任务也差不多, 我们需要将某位用户 (比如root) 全部文件的所有权更改成另一位用户 (比如Web服务器默认的Apache用户www-data), 那么就可以用 -user找出root拥有的所有文件, 然后用 -exec更改所有权。



你必须以超级用户的身份执行find命令才能够进行所有权的更改。

示例如下:

```
# find . -type f -user root -exec chown slynux {} \;
```

在这个命令中, {}是一个特殊的字符串, 与 -exec选项结合使用。对于每一个匹配的文件, {}会被替换成相应的文件名。例如, find命令找到两个文件test1.txt和test2.txt, 其所有者均为slynux, 那么find将会执行:

```
chown slynux {}
```

它会被解析为`chown slynux test1.txt`和`chown slynux test2.txt`。

另一个例子是将给定目录中的所有C程序文件拼接起来写入单个文件`all_c_files.txt`。我们可以用`find`找到所有的C文件，然后结合`-exec`使用`cat`命令：

```
$ find . -type f -name "*.c" -exec cat {} \;>all_c_files.txt
```

`-exec`之后可以接任何命令。`{}`表示一个匹配。对于任何一个匹配的文件名，`{}`会被该文件名所替换。

我们使用`>`操作符将来自`find`的数据重定向到`all_c_files.txt`文件，没有使用`>>`（追加）的原因是因为`find`命令的全部输出只是一个单数据流（`stdin`），而只有当多个数据流被追加到单个文件中的时候才有必要使用`>>`。

例如，用下列命令将10天前的`.txt`文件复制到`OLD`目录中：

```
$ find . -type f -mtime +10 -name "*.txt" -exec cp {} OLD \;
```

`find`命令同样可以采用类似的方法与其他命令结合起来。

-exec结合多个命令



我们无法在`-exec`参数中直接使用多个命令。它只能够接受单个命令，不过我们可以耍一个小花招。把多个命令写到一个shell脚本中（例如`command.sh`），然后在`-exec`中使用这个脚本：

```
-exec ./command.sh {} \;
```

`-exec`能够同`printf`结合起来生成有用的输出信息。例如：

```
$ find . -type f -name "*.txt" -exec printf "Text file: %s\n" {} \;
```

10. 让find跳过特定的目录

在搜索目录并执行某些操作的时候，有时为了提高性能，需要跳过一些子目录。例如，程序员会在版本控制系统（如Git）所管理的开发源码树中查找特定的文件，源代码层级结构总是会在每个子目录中包含一个`.git`目录（`.git`存储每个目录相关的版本控制信息）。因为与版本控制相关的目录对我们而言并没有什么用处，所以没必要去搜索这些目录。将某些文件或目录从搜索过程中排除的技术被称为“修剪”，其操作方法如下：

```
$ find devel/source_path \( -name ".git" -prune \) -o \( -type f -print \)
```

不使用`\(-type -print \)`，而是选择需要的过滤器

以上命令打印出不包括在`.git`目录中的所有文件的名称（路径）。

这里，`\(-name ".git" -prune \)`的作用是用于进行排除，它指明了`.git`目录应该被排除掉，而`\(-type f -print \)`指明了需要执行的动作。这些动作需要被放置在第二个语句块中（打印出所有文件的名称和路径）。

2.5 玩转 xargs

我们可以用管道将一个命令的stdout（标准输出）重定向到另一个命令的stdin（标准输入）。例如：

```
cat foo.txt | grep "test"
```

但是，有些命令只能以命令行参数的形式接受数据，而无法通过stdin接受数据流。在这种情况下，我们没法用管道来提供那些只有通过命令行参数才能提供的的数据。

那就只能另辟蹊径了。xargs是一个很有用的命令，它擅长将标准输入数据转换成命令行参数。xargs能够处理stdin并将其转换为特定命令的命令行参数。xargs也可以将单行或多行文本输入转换成其他格式，例如单行变多行或是多行变单行。

Bash黑客都喜欢单行命令。单行命令是一个命令序列，各命令之间不使用分号，而是使用管道操作符进行连接。精心编写的单行命令可以更高效、更简捷地完成任务。就文本处理而言，需要具备扎实的理论和实践才能够写出适合的单行命令解决方法。xargs就是构建单行命令的重要组成部分之一。

2.5.1 预备知识

xargs命令应该紧跟在管道操作符之后。它以标准输入作为主要的源数据流，并使用stdin并通过提供命令行参数来执行其他命令。例如：

```
command | xargs
```

2.5.2 实战演练

xargs命令把从stdin接收到的数据重新格式化，再将其作为参数提供给其他命令。

xargs可以作为一种替换方式，作用类似于find命令中的-exec参数。下面介绍一些借助xargs命令能够实现技巧。

❑ 将多行输入转换成单行输出

只需要将换行符移除，再用" "（空格）进行代替，就可以实现多行输入的转换。'\n'被解释成一个换行符，换行符其实就是多行文本之间的定界符。利用xargs，我们可以用空格替换掉换行符，这样一来，就能够将多行文本转换成单行文本：

```
$ cat example.txt # 样例文件
1 2 3 4 5 6
7 8 9 10
11 12
```

```
$ cat example.txt | xargs
1 2 3 4 5 6 7 8 9 10 11 12
```

❑ 将单行输入转换成多行输出

指定每行最大的参数数量n，我们可以将任何来自stdin的文本划分成多行，每行n个参数。

每一个参数都是由" " (空格)隔开的字符串。空格是默认的定义符。按照下面的方法可以将单行划分成多行:

```
$ cat example.txt | xargs -n 3
1 2 3
4 5 6
7 8 9
10 11 12
```

2.5.3 工作原理

`xargs`命令拥有数量众多且用法简单的选项,这使得它很适用于某些问题场合。让我们来看看如何能够巧妙地用这些选项来解决问题。

我们可以用自己的定义符来分隔参数。用 `-d` 选项为输入指定一个定制的定义符:

```
$ echo "splitXsplitXsplitXsplit" | xargs -d X
split split split split
```

在上面的代码中, `stdin` 是一个包含了多个 `X` 字符的字符串。我们可以用 `-d` 将 `X` 作为输入定义符。在这里,我们明确指定 `X` 作为输入定义符,而在默认情况下, `xargs` 采用内部字段分隔符 (IFS) 作为输入定义符。

同时结合 `-n`, 我们可以将输入划分成多行,而每行包含两个参数:

```
$ echo "splitXsplitXsplitXsplit" | xargs -d X -n 2
split split
split split
```

2.5.4 补充内容

我们已经从上面的例子中学到了如何将 `stdin` 格式化成为不同的输出形式以作为参数。现在让我们来学习如何将参数传递给命令。

1. 读取 `stdin`, 将格式化参数传递给命令

编写一个小型的定制版 `echo` 来更好地理解 `xargs` 提供命令行参数的方法:

```
#!/bin/bash
# 文件名: cecho.sh
```

```
echo $*'#'
```

当参数被传递给文件 `cecho.sh` 后, 它会将这些参数打印出来, 并以 `#` 字符作为结尾。例如:

```
$ ./cecho.sh arg1 arg2
arg1 arg2 #
```

让我们来看下面这个问题。

- 有一个包含着参数列表的文件 (每行一个参数)。我需要用两种方法将这些参数传递给一个命令 (比如 `cecho.sh`)。第一种方法, 需要每次提供一个参数:

```
./cecho.sh arg1
./cecho.sh arg2
./cecho.sh arg3
```

或者，每次需要提供两个或三个参数。提供两个参数时，可采用类似于下面这种形式的方法：

```
./cecho.sh arg1 arg2
./cecho.sh arg3
```

□ 第二种方法，需要一次性提供所有的命令参数：

```
./cecho.sh arg1 arg2 arg3
```

先别急着往下看，试着运行一下上面的命令，然后仔细观察输出结果。

上面的问题也可以用xargs来解决。我们有一个名为args.txt的参数列表文件，这个文件的内容如下：

```
$ cat args.txt
arg1
arg2
arg3
```

就第一个问题，我们可以将这个命令执行多次，每次使用一个参数：

```
$ cat args.txt | xargs -n 1 ./cecho.sh
arg1 #
arg2 #
arg3 #
```

每次执行需要X个参数的命令时，使用：

```
INPUT | xargs -n X
```

例如：

```
$ cat args.txt | xargs -n 2 ./cecho.sh
arg1 arg2 #
arg3 #
```

就第二个问题，我们可以执行这个命令，并一次性提供所有的参数：

```
$ cat args.txt | xargs ./ccat.sh
arg1 arg2 arg3 #
```

在上面的例子中，我们直接为特定的命令（例如cecho.sh）提供命令行参数。这些参数都只源于args.txt文件。但实际上除了它们外，我们还需要一些固定不变的命令参数。思考下面这种命令格式：

```
./cecho.sh -p arg1 -l
```

在上面的命令执行过程中，arg1是唯一的可变文本，其余部分都保持不变。我们应该从文件（args.txt）中读取参数，并按照下面的方式提供给命令：

```
./cecho.sh -p arg1 -l
./cecho.sh -p arg2 -l
./cecho.sh -p arg3 -l
```

xargs有一个选项-I, 可以提供上面这种形式的命令执行序列。我们可以用-I指定一个替换字符串, 这个字符串在xargs扩展时会被替换掉。当-I与xargs结合使用时, 对于每一个参数, 命令都会被执行一次。

试试下面的用法:

```
$ cat args.txt | xargs -I {} ./cecho.sh -p {} -l
-p arg1 -l #
-p arg2 -l #
-p arg3 -l #
```

-I {} 指定了替换字符串。对于每一个命令参数, 字符串{}会被从stdin读取到的参数所替换。使用-I的时候, 命令就似乎是在一个循环中执行一样。如果有三个参数, 那么命令就会连同{}一起被执行三次, 而{}在每一次执行中都会被替换为相应的参数。

2. 结合find使用xargs

xargs和find算是一对死党。两者结合使用可以让任务变得更轻松。不过, 人们通常却是一种错误的组合方式使用它们。例如:

```
$ find . -type f -name "*.txt" -print | xargs rm -f
```

这样做很危险。有时可能会删除不必要删除的文件。我们没法预测分隔find命令输出结果的定义符究竟是'\n'还是' '。很多文件名中都可能包含空格符, 而xargs很可能会误认为它们是定界符(例如, hell.txt.txt会被xargs误认为hell和text.txt)。

只要我们把find的输出作为xargs的输入, 就必须将-print0与find结合使用, 以字符null来分隔输出。

用find匹配并列出所有的.txt文件, 然后用xargs将这些文件删除:

```
$ find . -type f -name "*.txt" -print0 | xargs -0 rm -f
```

这样就可以删除所有的.txt文件。xargs-0将\0作为输入定界符。

3. 统计源代码目录中所有C程序文件的行数

统计所有C程序文件的行数是大多数程序员都会遇到的活儿。完成这项任务的代码如下:

```
$ find source_code_dir_path -type f -name "*.c" -print0 | xargs -0 wc -l
```

4. 结合stdin, 巧妙运用while语句和子shell

xargs只能以有限的几种方式来提供参数, 而且它也不能为多组命令提供参数。要执行一些包含来自标准输入的多个参数的命令, 可采用一种非常灵活的方法。我将这个方法称为子shell妙招(subshell hack)。一个包含while循环的子shell可以用来读取参数, 并通过一种巧妙的方式执行命令:

```
$ cat files.txt | ( while read arg; do cat $arg; done )
# 等同于 cat files.txt | xargs -I {} cat {}
```

在while循环中, 可以将cat \$arg替换成任意数量的命令, 这样我们就可以对同一个参数执行多项命令。我们也可以不借助管道, 将输出传递给其他命令。这个技巧能适用于各种问题环境。子shell内部的多个命令可作为一个整体来运行。

```
$ cmd0 | ( cmd1;cmd2;cmd3) | cmd4
```

如果cmd1是cd /, 那么就会改变子shell工作目录, 然而这种改变仅局限于子shell内部。cmd4则完全不知道工作目录发生了变化。

2.6 用 tr 进行转换

tr是UNIX命令行家工具箱中一件精美的小工具。它经常用来编写优美的单行命令, 作用不容小视。

tr可以对来自标准输入的字符进行替换、删除以及压缩。它可以将一组字符变成另一组字符, 因而通常也被称为转换 (translate) 命令。

2.6.1 预备知识

tr只能通过stdin (标准输入), 而无法通过命令行参数来接受输入。它的调用格式如下:

```
tr [options] set1 set2
```

将来自stdin的输入字符从set1映射到set2, 并将其输出写入stdout (标准输出)。set1和set2是字符类或字符集。如果两个字符集的长度不相等, 那么set2会不断重复其最后一个字符, 直到长度与set1相同。如果set2的长度大于set1, 那么在set2中超出set1长度的那部分字符则全部被忽略。

2.6.2 实战演练

将输入字符由大写转换成小写, 可以使用下面的命令:

```
$ echo "HELLO WHO IS THIS" | tr 'A-Z' 'a-z'
```

'A-Z' 和 'a-z' 都是集合。我们可以按照需要追加字符或字符类来构造自己定制的集合。

'ABD-}', 'aA.', 'a-ce-x' 以及 'a-c0-9' 等都是合法的集合。定义集合也很简单, 不需要去书写一长串连续的字符序列, 相反, 我们可以使用“起始字符-终止字符”这种格式。这种写法也可以和其他字符或字符类结合使用。如果“起始字符-终止字符”不是一个连续的字符序列, 那么它就会被视为一个包含了三个元素的集合 (“起始字符-终止字符”)。你可以使用像 '\t'、'\n' 这种特殊字符, 也可以使用其他ASCII字符。

2.6.3 工作原理

通过在tr中使用集合的概念, 我们可以轻松地将字符从一个集合映射到另一个集合中。让我们通过一则示例看看如何用tr进行数字加密和解密。

```
$ echo 12345 | tr '0-9' '9876543210'
87654 #已加密
```

```
$ echo 87654 | tr '9876543210' '0-9'
12345 #已解密
```

再看另外一个有趣的例子。

ROT13是一个著名的加密算法。在ROT13算法中,文本加密和解密都使用同一个函数。ROT13按照字母表排列顺序执行13个字母的转换。用tr进行ROT13加密:

```
$ echo "tr came, tr saw, tr conquered." | tr
'ABCDEFGHIJKLMNopQRSTUVWXYZabcdefghijklmnopqrstuvwxyz'
'NOPQRSTUVWXYZABCDEFGHIJKLMnopqrstuvwxyzabcdefghijklm'
```

得到输出:

```
ge pnxr, ge fnj, ge pbadhrerq.
```

对加密后的密文再次使用同样的ROT13函数,我们采用:

```
$ echo ge pnxr, ge fnj, ge pbadhrerq. | tr
'ABCDEFGHIJKLMNopQRSTUVWXYZabcdefghijklmnopqrstuvwxyz'
'NOPQRSTUVWXYZABCDEFGHIJKLMnopqrstuvwxyzabcdefghijklm.'
```

于是,得到输出:

```
tr came, tr saw, tr conquered.
```

tr还可以用来将制表符转换成空格:

```
$ cat text | tr '\t' ' '
```

2.6.4 补充内容

1. 用tr删除字符

tr有一个选项-d,可以通过指定需要被删除的字符集合,将出现在stdin中的特定字符清除掉:

```
$ cat file.txt | tr -d '[set1]'
#只使用set1,不使用set2
```

例如:

```
$ echo "Hello 123 world 456" | tr -d '0-9'
Hello world
# 将stdin中的数字删除并打印出来
```

2. 字符集补集

我们可以利用选项-c来使用set1的补集。-c [set]等同于定义了一个集合(补集),这个集合中的字符不包含在[set]中:

```
tr -c [set1] [set2]
```

set1的补集意味着这个集合中包含set1中没有的所有字符。

最典型的用法是从输入文本中将不在补集中的所有字符全部删除。例如:

```
$ echo hello 1 char 2 next 4 | tr -d -c '0-9 \n'
1 2 4
```

在这里，补集中包含了除数字、空格字符和换行符之外的所有字符。因为指定了-d，所以这些字符全部都被删除。

3. 用tr压缩字符

tr命令在很多文本处理环境中相当有用。多数情况下，连续的重复字符应该被压缩成单个字符，而经常需要进行的一项任务就是压缩空白字符。

tr的-s选项可以压缩输入中重复的字符，方法如下：

```
$ echo "GNU is not UNIX. Recursive right ?" | tr -s ' '
GNU is not UNIX. Recursive right ?
# tr -s '{set}'
```

让我们用一种巧妙的方式用tr将文件中的数字列表进行相加：

```
$ cat sum.txt
1
2
3
4
5
$ cat sum.txt | echo ${tr '\n' '+' } 0 }
15
```

这一招是如何起效的？

在上面的命令中，tr用来将'\n'替换成'+'，因此我们得到了字符串"1+2+3+...5+"，但是在字符串的尾部多了一个操作符+。为了抵消这个多出来的操作符，我们再追加一个0。

`\${operation}`执行算术运算，因此得到下面的字符串：

```
echo ${1+2+3+4+5+0 }
```

如果我们利用循环从文件中读取数字，然后再进行相加，那肯定得用好几行代码。如今我们只用一行就完成任务了。这种编写单行命令的技巧可通过实践才能习得。

4. 字符类

tr可以像使用集合一样使用各种不同的字符类，这些字符类如下所示。

- ☐ alnum: 字母和数字。
- ☐ alpha: 字母。
- ☐ cntrl: 控制（非打印）字符。
- ☐ digit: 数字。
- ☐ graph: 图形字符。
- ☐ lower: 小写字母。
- ☐ print: 可打印字符。
- ☐ punct: 标点符号。
- ☐ space: 空白字符。
- ☐ upper: 大写字母。
- ☐ xdigit: 十六进制字符。



可以按照下面的方式选择并使用所需的字符类：

```
tr [:class:] [:class:]
```

例如：

```
tr '[:lower:]' '[:upper:]'
```

2.7 校验和与核实

校验和 (checksum) 程序用来从文件中生成校验和密钥，然后利用这个校验和密钥核实文件的完整性。一份文件可以通过网络或任何存储介质分发到不同的地点。出于多种原因，数据有可能在传输过程中丢失了若干位，从而导致文件损坏。这种错误通常发生在从Internet上下载文件时，或者通过网络传输文件时，或者遭遇CD光盘损坏等。

因此，我们需要采用一些测试方法来确定接收到的文件是否存在错误。用于文件完整性测试的特定密钥就被称为校验和。

我们对原始文件和接收到的文件都进行校验和计算。通过比对两者的校验和，就能够核实接收到的文件是否正确。如果校验和（一个来自源位置的原始文件，另一个来自目的地的接收文件）相等，就意味着我们接收到了正确的文件，否则用户就不得不重新发送文件并再次比对校验和。

校验和对于编写备份脚本或系统维护脚本来说非常重要，因为它们都会涉及通过网络传输文件。通过使用校验和核实，我们就可以识别出那些在网络传输过程中出现损坏的文件，并重发这些文件，从而确保数据的完整性。

2.7.1 预备知识

最知名且使用最为广泛的校验和技术是md5sum和sha1sum。它们对文件内容使用相应的算法来生成校验和。下面就来看看如何从文件中生成校验和并核实文件的完整性。

2.7.2 实战演练

为了计算md5sum，使用下列命令：

```
$ md5sum filename
68b329da9893e34099c7d8ad5cb9c940 filename
```

如上所示，md5sum是一个32个字符的十六进制串。

将输出的校验和重定向到一个文件，然后用这个MD5文件核实数据的完整性：

```
$ md5sum filename > file_sum.md5
```

2.7.3 工作原理

md5sum校验和计算的方法如下：

```
$ md5sum file1 file2 file3 ..
```

当使用多个文件时，输出中会在每行中包含单个文件的校验和：

```
[checksum1] file1
[checksum1] file2
[checksum1] file3
```

可以按照下面的方法用生成的文件核实数据完整性：

```
$ md5sum -c file_sum.md5
# 这个命令会输出校验和是否匹配的消息
```

如果需要用所有的.md5信息来检查所有的文件，可以使用：

```
$ md5sum *.md5
```

与md5sum类似，SHA1是另一种常用的校验和算法。它从给定的输入文件中生成一个长度为40个字符的十六进制串。用来计算SHA1串的命令是shasum。其用法和md5sum的非常相似。只需要把先前讲过的那些命令中的md5sum替换成shasum就行了，另外，记住要将输入文件名从file_sum.md5改为file_sum.shal。

校验和对于核实下载文件的完整性非常有帮助。我们从Internet上下载的ISO镜像文件一般更容易出现错误。因而，为了检查接收文件正确与否，校验和得以广泛应用。校验和程序对同样的文件数据始终生成一模一样的校验和。

2.7.4 补充内容

对于多个文件，校验和同样可以发挥作用。现在就看看如何校验并核实多个文件。

对目录进行校验

校验和是从文件中计算得来的。对目录计算校验和意味着我们需要对目录中的所有文件以递归的方式进行计算。

它可以用命令md5deep或shaldeep来实现。首先，需要安装md5deep软件包以确保能找到这些命令。该命令的用法如下：

```
$ md5deep -rl directory_path > directory.md5
# -r 使用递归的方式
# -l 使用相对路径。默认情况下，md5deep会输出文件的绝对路径
```

或者，也可以结合find来递归计算校验和：

```
$ find directory_path -type f -print0 | xargs -0 md5sum >> directory.md5
用下面的命令进行核实：
$ md5sum -c directory.md5
```

2.8 排序、单一与重复

同文本文件打交道时，总离不开排序，那是因为对于文本处理任务而言，排序(sort)可以起到不小的作用。sort命令能够帮助我们对文本文件和stdin进行排序操作。通常，它会结合其他命令来生产所需要的输出。uniq是一个经常与sort一同使用的命令。它的作用是从文本或

stdin中提取单一的行。sort和uniq能够用来查找重复数据。这则攻略将演示sort和uniq命令的大多数用法。

2.8.1 预备知识

sort命令既可以从特定的文件，也可以从stdin中获取输入，并将输出写入stdout。uniq的工作模式和sort一样。

2.8.2 实战演练

我们可以按照下面的方式轻松地对一组文件（例如file1.txt和file2.txt）进行排序：

```
$ sort file1.txt file2.txt .. > sorted.txt
```

或是

```
$ sort file1.txt file2.txt .. -o sorted.txt
```

找出已排序文件中不重复的行：

```
$ cat sorted_file.txt | uniq> uniq_lines.txt
```

2.8.3 工作原理

sort和uniq可以派上用场的地方有很多。让我们来认识一些命令选项和使用方法。

按数字进行排序：

```
$ sort -n file.txt
```

按逆序进行排序：

```
$ sort -r file.txt
```

按月份进行排序（按照一月、二月、三月……这样的顺序）：

```
$ sort -M months.txt
```

还可以用下面的方法测试一个文件是否已经被排过序：

```
#!/bin/bash
#用途：排序
sort -C file ;
if [ $? -eq 0 ] ; then
    echo Sorted;
else
    echo Unsorted;
fi
```

#要检查是否按数字进行排序，应该使用sort -nC

如果需要合并两个排过序的文件，而且不需要对合并后的文件再进行排序，可以使用：

```
$ sort -m sorted1 sorted2
```

2.8.4 补充内容

1. 依据键或列进行排序

如果需要将下面的文本排序，我们可以按列来进行。

```
$ cat data.txt
1 mac      2000
2 winxp    4000
3bsd      1000
4 linux    1000
```

有很多方法可以对这段文本排序。目前它是按照序号（第一列）来排序的，我们也可以依据第二列和第三列来排序。

-k指定了排序应该按照哪一个键（key）来进行。键指的是列号，而列号就是执行排序时的依据。-r告诉sort命令按照逆序进行排序。例如：

```
# 依据第1列，以逆序形式排序
$ sort -nrk 1 data.txt
4 linux    1000
3bsd      1000
2 winxp    4000
1 mac      2000
# -nr表明按照数字，采用逆序形式排序

# 依据第2列进行排序

$ sort -k 2 data.txt
3bsd      1000
4 linux    1000
1 mac      2000
2 winxp    4000
```



留意用于按照数字顺序进行排序的选项-n，就依据字母表排序和依据数字顺序排序，sort命令对于字母表排序和数字排序有不同的处理方式。因此，如果要采用数字顺序排序，就应该明确地给出-n选项。

通常在默认情况下，键就是文本文件中的列。列与列之间用空格分隔。但有时候，我们需要使用特定范围内的一组字符（例如，key1=character4-character8）作为键。在这种情况下，必须明确地将键指定为某个范围的字符，这个范围可以用键起止的字符位置来表明。例如：

```
$ cat data.txt
1010hellothis
2189ababba
7464dfdfdfdf
$ sort -nk 2,3 data.txt
```

突出显示的字符将用作数值键。为了提取这个键，用字符的起止位置作为键的书写格式。用第一个字符作为键：

```
$ sort -nk 1,1 data.txt
```

为了使sort的输出与以\0作为参数终止符的xargs命令相兼容，采用下面的命令：

```
$ sort -z data.txt | xargs -0  
# 终止符\0使得xargs命令的使用更加安全
```

有时文本中可能会包含一些像空格之类的不必要的字符。如果需要忽略这些字符，并以字典序进行排序，可以使用：

```
$ sort -bd unsorted.txt
```

其中，选项 -b用于忽略文件中的前导空白字符，选项 -d用于指明以字典序进行排序。

2. uniq

uniq命令通过消除重复内容，从给定输入中（stdin或命令行参数文件）找出单行的行。它也可以用来找出输入中出现的重复行。uniq只能用于排过序的数据输入，因此，uniq要么使用管道，要么将排过序的文件作为输入，并总是以这种方式与sort命令结合起来使用。

你可以按照下面的方式从给定的输入数据中生成单行的行（所谓“单行的行”是指来自输入的所有行都会被打印出来，但是其中重复的行只会被打印一次）：

```
$ cat sorted.txt  
bash  
foss  
hack  
hack
```

```
$ uniq sorted.txt  
bash  
foss  
hack
```

或是

```
$ sort unsorted.txt | uniq
```

或是

```
$ sort -u unsorted.txt
```

只显示唯一的行（在输入文件中没有出现重复的行）：

```
$ uniq -u sorted.txt  
bash  
foss
```

或是

```
$ sort unsorted.txt | uniq -c
```

为了统计各行在文件中出现的次数，使用下面的命令：

```
$ sort unsorted.txt | uniq -c  
1 bash  
1 foss  
2 hack
```



找出文件中重复的行：

```
$ sort unsorted.txt | uniq -d
hack
```

我们可以结合 `-s` 和 `-w` 来指定键：

- `-s` 指定可以跳过前N个字符；
- `-w` 指定用于比较的最大字符数。

这个对比键被用作 `uniq` 操作的索引：

```
$ cat data.txt
u:01:gnu
d:04:linux
u:01:bash
u:01:hack
```

我们需要使用醒目的字符作为唯一的键。可以通过忽略前2个字符 (`-s 2`)，并使用 `-w` 选项 (`-w 2`) 指定用于比较的最大字符数的方式来选定该键。

```
$ sort data.txt | uniq -s 2 -w 2
d:04:linux
u:01:bash
```

我们将命令输出作为 `xargs` 命令的输入的时候，最好为输出的各行添加一个0值字节终止符。在将 `uniq` 命令的输入作为 `xargs` 的数据源时，同样应当如此。如果没有使用0值字节终止符，那么在默认情况下，`xargs` 命令会用空格作为定界符分割参数。例如，来自 `stdin` 的文本行 “this is a line” 会被 `xargs` 当做包含4个不同的参数，但实际上它只是一个单行而已。如果使用0值字节终止符，那么 `\0` 就被作为定界符，此时，包含空格的单行就能够被正确地解析为单个参数。

用 `uniq` 命令生成包含0值字节终止符的输出：

```
$ uniq -z file.txt
```

下面的命令将删除所有指定的文件，而这些文件的名字是从 `files.txt` 中读取的：

```
$ uniq -z file.txt | xargs -0 rm
```

如果某个文件名在文件中出现多次，`uniq` 命令只会将这个文件名写入 `stdout` 一次。

3. 用 `uniq` 生成字符串样式

这里有一个很有意思的问题：我们有一个包含重复字符的字符串，如何才能知道每个字符在字符串中出现的次数，并依照下面的格式输出字符串？

输入：ahebhaaa

输出：4a1b1e2h

任何一个字符只要出现重复，就在之前加上它在字符串中出现过的次数。我们可以结合运用 `uniq` 和 `sort` 来解决这个问题：

```
INPUT= "ahebhaaa"
OUTPUT=` echo $INPUT | sed 's/[^\n]/&\n/g' | sed '/^$/d' | sort | uniq
-c | tr -d '\n'`
echo $OUTPUT
```

将上面代码中的管道命令进行分解：

```
echo $INPUT # 将输入打印到stdout
sed 's/[^\.]/&\n/g'
```

在每个字符后追加一个换行符，使得每行只出现一个字符。这让我们可以用sort命令对字符进行排序。sort命令只能用于由换行符分隔的记录上。

- sed '/^\\$/d': 最后一个字符会被sed替换成“字符+\\n”，因此会多出一个换行符并在最后形成一个空行。而这个命令就用来删除这最后的空行。
- sort: 因为每行只有一个字符，所以可以进行排序，并将排序结果作为uniq的输入。
- uniq -c: 这个命令打印出每一行各重复了多少次（计数）。
- tr -d '\\n': 将输入中的空格和换行符删除，生成所要求的输出格式。

2.9 临时文件命名与随机数

编写shell脚本时，我们经常需要存储临时数据。最适合存储临时数据的位置是/tmp（该目录中的内容在系统重启后会被清空）。有两种方法可以为临时数据生成标准的文件名。

2.9.1 实战演练

tempfile命令只有在基于Debian的发布版中才能找到，如Ubuntu、Debian等。Linux的其他发布版中并没有这个命令。

下面的代码为变量temp_file赋值了一个临时文件名：

```
temp_file=$(tempfile)
```

用echo \$temp_file可以在终端中打印出这个临时文件名。

该输出看起来类似于/tmp/fileaZWm8Y。

有时候，我们可以用一个加带了随机数的文件名作为临时文件名。我们可以依照下面的方法来实现：

```
temp_file="/tmp/file-$RANDOM"
```

环境变量\$RANDOM总是返回一个随机数。

2.9.2 工作原理

就算不使用tempfile命令，我们也可以使用自己的临时文件名。多数有经验的UNIX程序员会使用下面的习惯用法：

```
temp_file="/tmp/var.$$"
```

.\$\$ 作为添加的后缀会被扩展成当前运行脚本的进程ID。

2.10 分割文件和数据

在某些情况下，必须把文件分割成多个更小的片段。早期像软盘之类的设备容量都很有限，那么，将文件分割得更小些，以便能够存储到多张磁盘中就显得至关重要了。不过如今我们分割文件就是出于其他目的了，比如为提高可读性、生成日志等。

2.10.1 工作原理

生成一个大小为100KB的测试文件 (data.file)：

```
$ dd if=/dev/zero bs=100k count=1 of=data.file
```

上面的命令会创建一个大小为100KB而文件内容全部是0的文件。

你可以指定分割大小，将文件分割成多个更小的文件：

```
$ split -b 10k data.file
$ ls
data.file xaa xab xac xad xae xaf xag xah xai xaj
```

上面的命令将data.file分割成多个文件，每一个文件大小为10KB。这些文件以xab、xac、xad这种方式命名。这意味着它们都有一个字母后缀。如果想以数字为后缀，可以另外使用-d参数。此外，使用-a length就可以指定后缀长度：

```
$ split -b 10k data.file -d -a 4
$ ls
data.file x0009 x0019 x0029 x0039 x0049 x0059 x0069 x0079
```

除了k (KB)，我们还可以使用M (MB)、G (GB)、c (byte)、w (word) 等。

2.10.2 补充内容

来看看split命令的其他选项。

为分割后的文件指定文件名后缀

上面那些分割后的文件都有一个文件名前缀“x”。我们也可以通过提供一个前缀名以使用我们自己的文件名前缀。split命令最后一个参数是PREFIX，其格式如下：

```
$ split [COMMAND_ARGS] PREFIX
```

我们加上分隔文件的前缀名来运行先前那个命令：

```
$ split -b 10k data.file -d -a 4 split_file
$ ls
data.file      split_file0002  split_file0005  split_file0008 strtok.c
split_file0000 split_file0003  split_file0006  split_file0009
split_file0001 split_file0004  split_file0007
```

如果不想按照数据块大小，而是需要根据行数来分割文件的话，可以使用-l no_of_lines；

```
$ split -l 10 data.file
# 分割成多个文件，每个文件包含10行
```

另一个有趣的工具是csplit。它能够依据指定的条件和字符串匹配选项对log文件进行分割。不妨看看这个工具是如何运作的。

csplit是split工具的一个变体。split只能根据数据大小或行数分割文件，而csplit可以根据文本自身的特点进行分割。是否存在某个单词或文本内容都可作为分割文件的条件。

看一个日志文件示例：

```
$ cat server.log
SERVER-1
[connection] 192.168.0.1 success
[connection] 192.168.0.2 failed
[disconnect] 192.168.0.3 pending
[connection] 192.168.0.4 success
SERVER-2
[connection] 192.168.0.1 failed
[connection] 192.168.0.2 failed
[disconnect] 192.168.0.3 success
[connection] 192.168.0.4 failed
SERVER-3
[connection] 192.168.0.1 pending
[connection] 192.168.0.2 pending
[disconnect] 192.168.0.3 pending
[connection] 192.168.0.4 failed
```

我们需要将这个日志文件分割成server1.log、server2.log和server3.log，这些文件的内容分别取自原文件中不同的SERVER部分。那么，可以使用下面的方法来实现：

```
$ csplit server.log /SERVER/ -n 2 -s {*} -f server -b "%02d.log" ; rm
server00.log
```

```
$ ls
server01.log server02.log server03.log server.log
```

有关这个命令的详细说明如下所示。

- ☐ /SERVER/ 用来匹配某一行，分割过程即从此处开始。
- ☐ /[REGEX]/ 表示文本样式。包括从当前行（第一行）直到（但不包括）包含“SERVER”的匹配行。
- ☐ {*} 表示根据匹配重复执行分割，直到文件末尾为止。可以用{参数}的形式来指定分割执行的次数。
- ☐ -s 使命令进入静默模式，不打印其他信息。
- ☐ -n 指定分割后的文件名后缀的数字个数，例如01、02、03等。
- ☐ -f 指定分割后的文件名前缀（在上面的例子中，server就是前缀）。
- ☐ -b 指定后缀格式。例如“%02d.log”，类似于C语言中printf的参数格式。在这里文件名=前缀+后缀=server + %02d.log。

因为分割后的第一个文件没有任何内容（匹配的单词就位于文件的第一行中），所以我们删除了server00.log。

2.11 根据扩展名切分文件名

有一些脚本是依据文件名进行各种处理的。我们可能会需要在保留扩展名的同时修改文件名、转换文件格式（保留文件名的同时修改扩展名）或提取部分文件名。shell具有的内建功能可以依据不同的情况来切分文件名。

2

2.11.1 实战演练

借助 % 操作符可以轻松将名称部分从“名称.扩展名”这种格式的文件名中提取出来。你可以按照下面的方法从sample.jpg中提取名称。

```
file_jpg='sample.jpg'
name=${file_jpg%.*}
echo File name is: $name
```

输出结果：

```
File name is: sample
```

下一个任务是将文件名的扩展名部分提取出来，这可以借助 # 操作符实现。

提取文件名中的 .jpg 并存储到变量file_jpg中：

```
extension=${file_jpg#*.}
echo Extension is: jpg
```

输出结果：

```
Extension is: jpg
```

2.11.2 工作原理

在第一个任务中，为了从“名称.扩展名”这种格式的文件名中提取名称，我们使用了%操作符。\${VAR%.*} 的含义是：

- ❑ 从 \$VARIABLE中删除位于 % 右侧的通配符（在前例中是.*）所匹配的字符串。通配符从右向左进行匹配。
- ❑ 给VAR赋值，VAR=sample.jpg。那么，通配符从右向左就会匹配到.jpg，因此，从 \$VAR 中删除匹配结果，就会得到输出“sample”。

% 属于非贪婪（non-greedy）操作。它从右到左找出匹配通配符的最短结果。还有另一个操作符%%，这个操作符与%相似，但行为模式却是贪婪的（greedy），这意味着它会匹配符合条件的最长的字符串。

例如，我们现在有这样一个文件：

```
VAR=hack.fun.book.txt
```

使用 % 操作符：

```
$ echo ${VAR%.*}
```

得到输出: `hack.fun.book`

操作符 `%` 使用 `*` 从右向左执行非贪婪匹配 (`.txt`)。

使用 `%%` 操作符:

```
$ echo ${VAR%*.}
```

得到输出: `hack`

操作符 `%%` 则用 `*` 从右向左执行贪婪匹配 (`.fun.book.txt`)。

在第二个任务中, 我们用 `#` 操作符从文件名中提取扩展名。这个操作符与 `%` 类似, 不过求值方向是从左向右。

`${VAR#*.}` 的含义是: 从 `$VARIABLE` 中删除位于 `#` 右侧的通配符 (即在前例中使用的 `*`。) 所匹配的字符串。通配符从左向右进行匹配。

和 `%%` 类似, `#` 也有一个相对应的贪婪操作符 `##`。

`##` 从左向右进行贪婪匹配, 并从指定变量中删除匹配结果。

来看一个例子:

```
VAR=hack.fun.book.txt
```

使用 `#` 操作符:

```
$ echo ${VAR#*.}
```

得到输出: `fun.book.txt`

操作符 `#` 用 `*` 从左向右执行非贪婪匹配 (`hack`)。

使用 `##` 操作符:

```
$ echo ${VAR##*.}
```

得到输出: `txt`

操作符 `##` 则用 `*` 从左向右执行贪婪匹配 (`hack.fun.book`)。



因为文件名中可能包含多个 `'.'` 字符, 所以相较于 `#`, `##` 更适合于从文件名中提取扩展名。 `##` 执行的是贪婪匹配, 因而总是能够准确地提取出扩展名。

这里有个能够提取域名不同部分的实用案例。假定 `URL="www.google.com"`:

```
$ echo ${URL%*.} # 移除 * 所匹配的最右边的内容
www.google
```

```
$ echo ${URL%%*.} # 将从右边开始一直匹配到最左边的 *. 移除 (贪婪操作符)
www
```

```
$ echo ${URL#*.} # 移除 * 所匹配的最左边的内容
google.com
```

```
$ echo ${URL##*.} # 将从左边开始一直匹配到最右边的 *. 移除 (贪婪操作符)
com
```

2.12 批量重命名和移动

重命名多个文件是我们经常会碰到的一项工作。举个简单的例子，当你把照片从数码相机传输到你的计算机之后，你可能会删除其中某些不如意的部分，这会使图像文件的编号变得不再连续。于是你会想使用特定的前缀和连续的数字对它们进行重命名。我们当然可以借助第三方软件执行这类重命名操作，但是我们也可以使用Bash命令在短短几秒钟之内完成同样的工作。

另一件经常要做的工作是，将文件名中包含某个特定部分（例如相同的前缀）或者具有特定类型的所有文件移动到指定的文件夹中。让我们看看如何用脚本来执行这些操作。

2.12.1 预备知识

rename命令利用Perl正则表达式修改文件名。综合运用find、rename和mv，我们其实可以完成很多操作。

2.12.2 实战演练

用特定的格式重命名当前目录下的图像文件，最简单的方法是使用下面的脚本：

```
#!/bin/bash
#文件名 rename.sh
#用途：重命名.jpg和.png文件

count=1;
for img in *.jpg *.png
do
new=image-$count.${img##*.}

mv "$img" "$new" 2> /dev/null

if [ $? -eq 0 ];
then

echo "Renaming $img to $new"
let count++

fi

done
```

输出如下：

```
$ ./rename.sh
Renaming hack.jpg to image-1.jpg
Renaming new.jpg to image-2.jpg
Renaming next.jpg to image-3.jpg
```

该脚本将当前目录下所有的.jpg和.png文件重命名，新文件名的格式为image-1.jpg、image-2.jpg、image-3.jpg、image-4.png等，依次类推。

2.12.3 工作原理

在前文列出的重命名脚本中，我们使用for循环对所有扩展名为.jpg的文件进行迭代，利用通配符*.jpg和*.png匹配所有的JPEG和PNG文件。我们可以稍稍改进一下扩展名匹配。.jpg只能匹配小写的扩展名，而将.jpg改成.[jJ][pP][gG]就会使这项匹配不区分大小写，这时，它不仅可以匹配扩展名为jpg的文件，也可以匹配扩展名为JPG或Jpg的文件。在Bash中，字符可以包含在[]之中，这意味着只要某个字符属于[]中的那个字符集合，[]就能够匹配它。

上面代码中的for img in *.jpg *.png会自动扩展为：

```
for img in hack.jpg new.jpg next.jpg
```

为了跟踪图像编号，我们初始化变量count=1。下一步就是用mv命令重命名文件。因此需要构造出新的文件名。\${img##*.}对处于当前循环中的文件名进行解析并获得文件扩展名（请参看2.11.2节中对于\${img##*.}的解释）。

let count++ 用来在每次循环中增加文件编号。

你可以看到2>操作符将mv命令的标准错误(stderr)重定向到/dev/null。这是为了防止错误信息被打印到终端。

如果通配符*.png和*.jpg没有匹配到任何图像文件，那么shell会将通配符解释成一个字符串。从上面的输出中可以看到并不存在扩展名为.png的文件。因此，如果将*.png作为一个文件名并执行mv *.png image-X.png，肯定会产生错误。if [\$? -eq 0]语句用来检查退出状态(\$?)。如果最近执行的命令没有错误，那么\$?的值是0，否则它会返回非0值。当mv命令出错时，也会返回非0值，因此不会出现提示信息“Renaming file”，计数也不会增加。

还有许多其他的执行重命名操作的方法，这里再举几个简单的例子。

将*.JPG更名为*.jpg:

```
$ rename *.JPG *.jpg
```

将文件名中的空格替换成字符“_”:

```
$ rename 's/ /_/g' *
```

's/ /_/g'用于替换文件名，而*是用于匹配目标文件的通配符，它也可以以*.txt或其他样式出现。

另外，你也可以转换文件名的大小写:

```
$ rename 'y/A-Z/a-z/' *
```

```
$ rename 'y/a-z/A-Z/' *
```

为了将所有的.mp3文件移入给定的目录，可以使用:

```
$ find path -type f -name "*.mp3" -exec mv {} target_dir \;
```

将所有文件名中的空格替换为字符“_”:

```
$ find path -type f -exec rename 's/ /_/g' {} \;
```

2.13 拼写检查与词典操作

Linux大多数的发行版都含有一份词典文件。然而，我发现几乎没人在意过这个文件，更别提使用它了。另外，有一个被称为`aspell`的命令行工具，其作用是进行拼写检查。让我们通过一些脚本，看看如何使用词典文件和拼写检查工具。

2

2.13.1 实战演练

目录 `/usr/share/dict/` 包含了一些词典文件。“词典文件”就是包含了一个词典单词列表的文本文件。我们可以利用这个列表来检查某个单词是否为词典中的单词。

```
$ ls /usr/share/dict/  
american-english british-english
```

为了检查给定的单词是否为词典中的单词，可以使用下面的脚本：

```
#!/bin/bash  
#文件名: checkword.sh  
word=$1  
grep "^$1$" /usr/share/dict/british-english -q  
if [ $? -eq 0 ]; then  
    echo $word is a dictionary word;  
else  
    echo $word is not a dictionary word;  
fi
```

这个脚本的用法如下：

```
$ ./checkword.sh ful  
ful is not a dictionary word  
  
$ ./checkword.sh fool  
fool is a dictionary word
```

2.13.2 工作原理

在`grep`中，`^`标记着单词的开始，`$`标记着单词的结束。

`-q` 禁止产生任何输出。

另外，我们也可以使用拼写检查命令`aspell`来核查某个单词是否在词典中：

```
#!/bin/bash  
#文件名: aspellcheck.sh  
word=$1  
output=`echo \"$word\" | aspell list`  
if [ -z $output ]; then  
    echo $word is a dictionary word;  
else  
    echo $word is not a dictionary word;  
fi
```

当给定的输入不是一个词典单词时，`aspell list`命令产生输出文本，反之则不产生任何

输出。-z用于确认 \$outpu 是否为空。

列出文件中以特定单词起头的所有单词：

```
$ look word filepath
```

或者使用

```
$ grep "^word" filepath
```

在默认情况下，如果没有给出文件参数，look命令会使用默认词典 (/usr/share/dict/words) 并返回输出。

```
$ look word
# 像这样使用时，look命令以默认词典作为文件参数
```

例如：

```
$ look android
android
android's
androids
```

2.14 交互输入自动化

就编写自动化工具或测试工具而言，实现命令的交互输入自动化极其重要。在很多情况下，我们要同一些交互式读取输入的命令打交道。“交互式输入”是指只有当命令要求获取输入时，才由用户手动键入。下面的例子就是一个要求提供交互式输入的命令执行过程。

```
$ command
Enter a number: 1
Enter name : hello
You have entered 1,hello
```

2.14.1 预备知识

能够自动接受上例中那种交互式输入的自动化工具，对于本地命令或远程应用来说非常有用。让我们看看如何实现自动化。

2.14.2 实战演练

试着思考交互式输入的过程。参照之前的代码，我们可以将这些步骤描述如下：

```
1[Return]hello[Return]
```

观察实际通过键盘输入的字符，可以将上面的1、Return、hello以及Return转换为以下字符串：

```
"1\nhello\n"
```

我们按下回车键时会发送 \n。添加上 \n后，就得到了发送给stdin的实际字符串。

因此，通过发送与用户输入等值的字符串，我们就可以实现在交互式进程中自动传递输入。

2.14.3 工作原理

先写一个读取交互式输入脚本，然后用这个脚本进行自动化的演示：

```
#!/bin/bash
#文件名: interactive.sh
read -p "Enter number:" no;
read -p "Enter name:" name
echo You have entered $no, $name;
```

按照下面的方法向命令自动发送输入：

```
$ echo -e "1\nhello\n" | ./interactive.sh
You have entered 1, hello
```

看来我们自己制作的输入生效了。

我们用echo -e来生成输入序列。如果输入内容比较多，那么可以用单独的输入文件结合重定向操作符来提供输入。

```
$ echo -e "1\nhello\n" > input.data
$ cat input.data
1
hello
```

制作输入文件时，你也可以不用echo命令：

```
$ ./interactive.sh < input.data
```

这个方法是从文件中导入交互式输入数据。

如果你是一位逆向工程师，那你可能同缓冲区溢出攻击打过交道。要实施攻击，我们需要将十六进制形式的shellcode（例如“\xeb\x1a\x5e\x31\xc0\x88\x46”）进行重定向。这些字符没法直接通过键盘输入，因为键盘上并没有其对应的键值。因此，我们应该使用：

```
echo -e "\xeb\x1a\x5e\x31\xc0\x88\x46"
```

这条命令会将shellcode重定向到有缺陷的可执行文件中。

我们已经描述了一种方法，它通过stdin将所需的文本进行重定向，从而实现交互式输入程序自动化。但是我们并没有检查所发送的输入内容。我们期望程序以特定（固定）的次序处理我们所发送的输入。如果程序对于输入采取随机或其他处理次序，或者甚至不要求输入某些内容，那么之前的方法就要出问题了。它会发送不符合程序要求的错误输入。为了处理动态输入并通过检查程序运行时的输入需求来提供输入内容，我们要使用一个出色的工具expect。expect命令可以根据适合的输入要求提供适合的输入。让我们看看如何使用expect。

2.14.4 补充内容

交互式输入自动化也可以用其他方法实现。expect脚本就是其中之一。下面来认识一下这种方法。

用expect实现自动化

在默认情况下，expect并没有附带于多数常见的Linux发行版中。expect必须要用软件包管理器手动安装。

expect等待特定的输入提示，通过检查输入提示来发送数据。

```
#!/usr/bin/expect
#文件名: automate_expect.sh
spawn ./interactive .sh
expect "Enter number:"
send "1\n"
expect "Enter name:"
send "hello\n"
expect eof
```

运行方法：

```
$ ./automate_expect.sh
```

在这个脚本中：

- spawn参数指定需要自动化哪一个命令；
- expect参数提供需要等待的消息；
- send是要发送的消息；
- expect eof指明命令交互结束。



本章内容

- 生成任意大小的文件
- 文本文件的交集与差集
- 查找并删除重复文件
- 创建长路径目录
- 文件权限、所有权和粘滞位
- 创建不可修改文件
- 批量生成空白文件
- 查找符号链接及其指向目标
- 列举文件类型统计信息
- 环回文件与挂载
- 生成ISO文件及混合ISO
- 查找文件差异并进行修补
- head与tail —— 打印文件的前10行或后10行
- 只列出目录的其他方法
- 在命令行中用pushd和popd快速定位
- 统计文件的行数、单词数和字符数
- 打印目录树

3.1 简介

UNIX将操作系统中的一切都视为文件。文件与每一个操作息息相关，而我们可以利用它们进行各种与系统或进程相关的处理工作。例如，我们所使用的命令终端就是和一个设备文件关联在一起的。我们可以通过写入特定终端所对应的设备文件来实现向终端写入信息。有各种不同形式的文件，比如目录、普通文件、块设备、字符设备、符号链接、套接字和命名管道等。文件名、大小、文件类型、文件内容修改时间（modification time）、文件访问时间（access time）、文件属性更改时间（change time）、i节点、链接以及文件所在的文件系统等都是文件的属性。本章主要包含了多个文件相关操作及属性的实战攻略。

3.2 生成任意大小的文件

由于各种可能的原因，你也许需要生成一个包含随机数据的文件。这可能是用于执行测试的测试文件，比如用一个大文件作为输入来测试应用程序的效率，也可能是测试文件分割，或是创建环回文件系统（环回文件自身可以包含文件系统，这种文件可以像物理设备一样使用mount命令进行挂载）。专门写一个程序来创建这些文件可是件麻烦事，所以我们用一些通用工具帮忙解决。

实战演练

创建特定大小的大文件最简单的方法就是利用dd命令。dd命令会克隆给定的输入内容，然后将一模一样的副本写入到输出。stdin、设备文件、普通文件等都可作为输入，stdout、设备文件、普通文件等也可作为输出。下面是使用dd命令的一个示例：

```
$ dd if=/dev/zero of=junk.data bs=1M count=1
1+0 records in
1+0 records out
1048576 bytes (1.0 MB) copied, 0.00767266 s, 137 MB/s
```

该命令会创建一个1MB大小的文件junk.data。来看一下命令参数：if代表输入文件（input file），of代表输出文件（output file），bs代表以字节为单位的块大小（block size, BS），count代表需要被复制的块数。

我们将bs指定为1MB，count指定为1，于是得到了一个大小为1MB的文件。如果把bs设为2MB，count设为2，那么总文件大小就是4MB。

块大小可以使用各种计量单位。表3-1中任意一个字符都可以置于数字之后来指定字节数。

表 3-1

单元大小	代 码
字节 (1B)	c
字 (2B)	w
块 (512B)	b
千字节 (1024B)	k
兆字节 (1024KB)	M
吉字节 (1024MB)	G

按照这种方法，我们就可以生成任意大小的文件。包括MB在内的表3-1中给出的计量单位都可以使用。

/dev/zero是一个字符设备，它会不断返回0值字节（\0）。

如果不指定输入参数（if），默认情况下dd会从stdin中读取输入。与之类似，如果不指定输出参数（of），则dd会将stdout作为默认输出。

也可以用dd命令传输大量数据并观察命令输出来测量内存的操作速度（例如，在前一个例子中所示的1048576 bytes (1.0 MB) copied, 0.00767266 s, 137 MB/s）。

3.3 文本文件的交集与差集

交集（intersection）和差集（set difference）操作在集合论相关的数学课上经常会被用到。不过，对文本进行类似的操作在某些情况下也很有用。

3.3.1 预备知识

`comm`命令可用于两个文件之间的比较。它有很多不错的选项可用来调整输出，以便我们执行交集、求差（difference）以及差集操作^①。

- 交集：打印出两个文件所共有的行。
- 求差：打印出指定文件所包含的且不相同的那些行。
- 差集：打印出包含在文件A中，但不包含在其他指定文件中的那些行。

3.3.2 实战演练

需要注意的是`comm`必须使用排过序的文件作为输入。请看看下面的例子：

```
$ cat A.txt
apple
orange
gold
silver
steel
iron

$ cat B.txt
orange
gold
cookies
carrot

$ sort A.txt -o A.txt ; sort B.txt -o B.txt
$ comm A.txt B.txt
apple          carrot
               cookies
               gold
iron           orange
silver
steel
```

输出的第一列包含只在A.txt中出现的行，第二列包含只在B.txt中出现的行，第三列包含A.txt和B.txt中相同的行。各列以制表符（\t）作为定界符。

有一些选项可以按照我们的需求进行格式化输出，例如：

- -1 从输出中删除第一列；
- -2 从输出中删除第二列；
- -3 从输出中删除第三列。

^① 假设现在有两个文件A和B，内容分别是：A(1,2,3)，B(3,4,5)。那么，对这两个文件进行操作的结果如下。

交集：3
求差：1,2,4,5
差集（A）：1,2。

为了打印两个文件的交集，我们需要删除第一列和第二列，只打印出第三列：

```
$ comm A.txt B.txt -1 -2
gold
orange
```

打印出两个文件中不相同的行：

```
$ comm A.txt B.txt -3
apple
        carrot
        cookies
iron
silver
steel
```

在comm命令中用-3删除第三列，输出第一列和第二列。第一列包含那些只出现在A.txt中的行，第二列包含那些只出现在B.txt中的行。由于输出中只包含了两列，它在此处的用处也不大。对于每个唯一的行，在列中都会对应一个空白字段。因此在同一行上，两列的内容不会重复。为了获得一个可用的输出文本格式，我们需要将空白字段删除，把两列合成一列：

```
apple
carrot
cookies
iron
silver
steel
```

要生成这样的输出，就得删除行首的\t字符。可以用下面的方法删除\t，从而合并成单列：

```
$ comm A.txt B.txt -3 | sed 's/^\\t//'
```

```
apple
carrot
cookies
iron
silver
steel
```

sed命令通过管道获取comm的输出。它删除行首的\t字符。sed中的s表示替换(substitute)。
/^\\t/ 匹配行前的\t(^是行首标记)。
// (两个/操作符之间没有任何字符)是用来替换行首的\t的字符串。如此一来，就删除了所有行首的\t。

接下来讲解两个文件的差集操作。

差集操作允许你比较两个文件，并打印出只在A.txt或B.txt中出现的行。当A.txt和B.txt作为comm命令的参数时，输出中的第一列是A.txt相对于B.txt的差集，第二列是B.txt相对于A.txt的差集。

通过删除不需要的列，我们就可以分别得到A.txt和B.txt的差集。

□ A.txt的差集

```
$ comm A.txt B.txt -2 -3
```

-2 -3 删除第二列和第三列。

□ B.txt的差集

```
$ comm A.txt B.txt -1 -3
-1 -3 删除第一列和第三列。
```

3.4 查找并删除重复文件

重复文件是同一个文件的多个副本。有时候，我们需要删除重复的文件，只保留单个副本。通过查看文件内容来识别重复文件是件挺有意思的任务。可以结合多种shell工具来完成这项任务。在这则攻略中，我们讨论的是查找重复文件并基于查找结果执行相关的操作。

3.4.1 预备知识

重复文件指的是那些虽然名字不同但内容却一模一样的文件。我们可以通过比较文件内容来识别它们。校验和是依据文件内容来计算的，内容相同的文件自然会生成相同的校验和，因此，我们可以通过比较校验和来删除重复文件。

3.4.2 实战演练

按照下面的方法创建一些测试文件：

```
$ echo "hello" > test ; cp test test_copy1 ; cp test test_copy2;
$ echo "next" > other;
# test_copy1和test_copy2都是test文件的副本
```

删除重复文件的脚本代码如下：

```
#!/bin/bash
#文件名: remove_duplicates.sh
#用途: 查找并删除重复文件，每一个文件只保留一个样本

ls -ls | awk 'BEGIN {
getline;getline;
name1=$8; size=$5
}
{ name2=$8;
if (size==$5)
{
md5sum "name1 | getline; csum1=$1;
md5sum "name2 | getline; csum2=$1;
if ( csum1==csum2 )
{print name1; print name2 }
};

size=$5; name1=name2;
}' | sort -u > duplicate_files
```



```
cat duplicate_files | xargs -I { } md5sum { } | sort | uniq -w 32 | awk
'{ print "'$2"$" }' | sort -u > duplicate_sample

echo Removing..
comm duplicate_files duplicate_sample -2 -3 | tee /dev/stderr | xargs
rm
echo Removed duplicates files successfully.

执行方式:

$ ./remove_duplicates.sh
```

3.4.3 工作原理

前文中的shell脚本会找出某个目录中同一文件的所有副本，然后保留单个副本的同时删除其他副本。让我们研究一下这个脚本的工作原理。ls -ls对当前目录下的所有文件按照文件大小进行排序，并列出文件的详细信息。awk读取ls -ls的输出，对行列进行比较，找出重复文件。

下面是代码的执行逻辑。

- 我们将文件依据大小排序并列出，这样大小接近的文件就会排列在一起。识别大小相同的文件是我们查找重复文件的第一步。接下来，计算这些文件的校验和。如果校验和相同，那么这些文件就是重复文件，将被删除。
- 在从文件中读取文本行之前，首先要执行awk的BEGIN{}语句块。读取文本行的工作在{}语句块中进行，读取并处理完所有的行之后，执行END{}语句块。ls -ls的输出如下：

```
total 16
4 -rw-r--r-- 1 slynux slynux 5 2010-06-29 11:50 other
4 -rw-r--r-- 1 slynux slynux 6 2010-06-29 11:50 test
4 -rw-r--r-- 1 slynux slynux 6 2010-06-29 11:50 test_copy1
4 -rw-r--r-- 1 slynux slynux 6 2010-06-29 11:50 test_copy2
```

- 第1行输出告诉我们文件数量，这个信息在本例中没什么用处。我们用getline读取第1行，然后丢弃。由于需要对每一行及其下一行来进行文件大小比对，因此用getline读取长文件列表的第一行，并存储文件名和大小（它们分别是第8列和第5列）。这样我们就提前读取到了一行。接下来，awk进入{}语句块（在这个语句块中读取余下的文本行），对于读取到的每一行文本，都会执行该语句块。它将从当前行中读取到的文件大小与之前存储在变量size中的值进行比较。如果相等，那就意味着两个文件至少在大小上是相同的，随后再用md5sum执行进一步的检查。

至此，我们已试着采用了一些技巧来逐步解决问题。

在awk中，外部命令的输出可以用下面的方法读取：

```
*cmd* | getline
```

随后，我们就可以在\$0中获取命令的输出，在\$1,\$2,...\$n中获取命令输出中的每一列。我们将文件的md5sum保存在变量csum1和csum2中。变量name1和name2保存文件列表中位置连续的文件名。如果两个文件的校验和相同，那它们肯定是重复文件，而它们的文件名会被打印出来。

我们需要从每组重复文件中找出一个文件，使得在保留一个副本的同时，能够删除其他重复的文件。我们计算重复文件的md5sum，从每一组重复文件中打印出一个文件。这是通过-w 32比较每一行的md5sum（md5sum输出中的前32个字符，md5sum的输出通常由32个字符的散列值和文件名组成），然后找出那些不相同的行。因此，每组重复文件中的一个采样就被写入duplicate_sample。

现在需要将duplicate_files中列出的且未被列于duplicate_sample之内的全部文件删除。comm命令打印出未在duplicate_sample中列出且被duplicate_files包含的文件。

因此，我们可以使用差集操作（请参考3.3节）。

comm通常只接受排序过的文件。所以，在重定向到duplicate_files和duplicate_sample之前，首先用sort -u作为一个过滤器。

tee命令在这里有一个妙用：它在将文件名传递给rm命令的同时，也起到了print的功能。tee将来自stdin的行写入文件，并将其发送到stdout。我们也可以将文本重定向到stderr来实现终端打印功能。/dev/stderr是对应于stderr（标准错误）的设备。通过重定向到stderr设备文件，来自stdin的文本将会以标准错误的形式出现在终端中。

3

3.4.4 参考

- 4.7节讲解了awk命令。
- 2.7节讲解了md5sum命令。

3.5 创建长路径目录

有时候，我们需要创建一个空目录树。如果给定路径中包含目录，那么还必须检查这些目录是否已经存在。这会使得代码变得臃肿而低效。这则攻略就是来解决这个问题的。

3.5.1 预备知识

mkdir命令用于创建目录。例如：

```
$ mkdir dirpath
```

如果目录已经存在，会返回“File exists”错误信息：

```
mkdir: cannot create directory 'dir_name': File exists
```

如果给你一个目录路径（/home/slynx/test/hello/child），其中/home/slynx已经存在，那么我们需要创建路径中余下的目录（/home/slynx/test、/home/slynx/test/hello和/home/slynx/test/hello/child）。

下面的代码可找出路径中的每个目录是否存在：

```
if [ -e /home/slynx ]; then
# 创建下一级目录
fi
```

-e是一个用在条件判断[]中的参数，可用以判断某个文件是否存在。在类UNIX系统中，目录同样是一种文件。如果该文件存在，[-e FILE_PATH]返回真。

3.5.2 实战演练

下面的命令序列可用于创建多级目录树：

```
$ mkdir /home 2> /dev/null
$ mkdir /home/slynux 2> /dev/null
$ mkdir /home/slynux/test 2> /dev/null
$ mkdir /home/slynux/test/hello 2> /dev/null
$ mkdir /home/slynux/test/hello/child 2> /dev/null
```

如果遇到“Directory exists”这种错误，该命令会被忽略，错误信息通过2>被重定向到/dev/null。不过这种方法用起来太烦琐，而且也不是标准做法。进行这一操作的标准单行命令如下：

```
$ mkdir -p /home/slynux/test/hello/child
```

这条单行命令足以取代之前的那5条命令。它会忽略所有已存在的目录，同时创建缺失的部分。

3.6 文件权限、所有权和粘滞位

文件权限和所有权是UNIX/Linux文件系统（如ext文件系统）最显著的特性之一。在UNIX/Linux平台工作时，我们经常会碰到与文件权限及所有权相关的问题。这则攻略就考察了文件权限和所有权的不同用例。

3.6.1 预备知识

Linux系统中的每一个文件都与多种权限类型相关联。在这些权限中，我们通常要和三类权限打交道（用户、用户组以及其他实体）。

用户（user）是文件的所有者。用户组（group）是多个用户的集合，系统允许这些用户对文件进行某些形式的访问。其他用户（others）是除用户或用户组之外的任何用户。

用命令ls -l可以列出文件的权限：

```
-rw-r--r-- 1 slynux slynux 2497 2010-02-28 11:22 bot.py
-rw-r--r-- 1 slynux slynux 16237 2010-02-06 21:42 c9.php
drwxr-xr-x 2 slynux slynux 4096 2010-05-27 14:31a.py
-rw-r--r-- 1 slynux slynux 539 2010-02-10 09:11 cl.pl
```

第1列输出明确了后面的输出。其中第一个字母的对应关系如下所示。

- ☐ “-” —— 普通文件。
- ☐ “d” —— 目录。
- ☐ “c” —— 字符设备。
- ☐ “b” —— 块设备。
- ☐ “l” —— 符号链接。
- ☐ “s” —— 套接字。
- ☐ “p” —— 管道。

剩下的部分可以划分成三组，每组3个字符（———）。第一组的3个字符（——）对应用户权限（所有者），第二组对应用户组权限，第三组对应其他用户权限。这9个字符（即9个权限）中的每一个字符指明是否其设置了某种权限。如果设置了权限，对应位置上会出现一个字符，否则就以一个“-”表明没有设置对应的权限。

让我们来看一下每个字符组对用户、用户组以及其他用户的含义。

用户

权限序列：rwx-----

第一个字符指定用户是否拥有文件的读权限。如果为用户设置了读权限，r将出现在第一个字符的位置上。第二个字符指定了写（修改）权限（w），第三个字符指定了用户是否拥有执行权限（x，即运行该文件的权限）。可执行文件通常会设置执行权限。用户还有一个称为setuid（s）的特殊权限，它出现在执行权限（x）的位置。setuid权限允许用户以其拥有者的权限来执行可执行文件，即使这个可执行文件是由其他用户运行的。

具有setuid权限的文件的权限序列如下：

-rws-----

目录同样也有读、写、执行权限。不过对于目录来说，读、写、执行权限的含义有点不一样。

- 目录的读权限（r）允许读取目录中文件和子目录的列表。
- 目录的写权限（w）允许在目录中创建或删除文件或目录。
- 目录的执行权限（x）指明是否可以访问目录中的文件和子目录。

用户组

权限序列：---rwx---

第二组字符指定了组权限。组权限的rwx的含义和用户权限中的一样。组权限并没有setuid，但是有一个setgid位。它允许以同该目录拥有者所在组相同的有效组权限来允许可执行文件。但是这个组 and 实际发起命令的用户组未必相同。组权限的权限序列如下：

----rws---

其他用户

权限序列：-----rwx

最后三个字符是其他用户权限。和用户以及用户组一样，其他用户也有读、写、执行权限，但是并没有S权限（如setuid和setgid）。

目录有一个特殊的权限，叫做粘滞位（sticky bit）。当一个目录设置了粘滞位，只有创建该目录的用户才能删除目录中的文件，即使用户组和其他用户也有写权限。粘滞位出现在其他用户权限中的执行权限（x）位置。它使用t或T来表示。如果没有设置执行权限，但设置了粘滞位，那么使用t；如果同时设置了执行权限和粘滞位，就使用T。

例如：

-----rwt , -----rwt

默认设置目录粘滞位的典型例子就是 /tmp。粘滞位属于一种写保护。

在ls -l的每一行输出中，字符串slynux slynux对应所属用户和所属用户组。第一个slynux代表所属用户，第二个slynux代表用户组的所有者。

3.6.2 实战演练

为了设置文件权限，可使用chmod命令。

假设需要设置权限：rwx rw- r--

就可以使用chmod按照下面的方法设置：

```
$ chmod u=rwx g=rw o=r filename
```

在这里：

□ u = 指定用户权限

□ g = 指定用户组权限

□ o = 指定其他实体权限

如果需要给文件添加权限，可以对用户、用户组和其他用户用+进行添加，用-删除权限。一个文件已经具有权限rwx rw- r--，现在需要增加可执行权限，方法如下：

```
$ chmod o+x filename
```

给所有权限类别（即用户、用户组和其他用户）增加可执行权限：

```
$ chmod a+x filename
```

其中，a表示全部（all）。

如果需要删除权限，则使用-，例如：

```
$ chmod a-x filename
```

也可以用八进制数来设置权限。权限由3位八进制数来表示，每一位按顺序分别对应用户、用户组和其他用户。

读、写和执行权限都有与之对应的唯一的八进制数：

□ r-- = 4

□ -w- = 2

□ --x = 1

我们可以将权限序列的八进制值相加来获得所需的权限组合，例如：

□ rwx = 4 + 2 + 6

□ r-x = 4 + 1 = 5

权限序列 rwx rw- r--的数字表示形式如下：

□ rwx = 4 + 2 + 1 = 7

□ rw- = 4 + 2 = 6

□ r-- = 4

因此，rwx rw- r-- 等于764，那么使用八进制值设置权限的命令为：

```
$ chmod 764 filename
```

3.6.3 补充内容

让我们再看一些其他有关文件和目录的操作任务。

1. 更改所有权

要更改文件的所有权，可以使用chown命令：

```
$ chown user.group filename
```

例如：

```
$ chown slynux.slynux test.sh
```

在这里，slynux既是用户名，也是用户组名。

2. 设置粘滞位

粘滞位是一种应用于目录的权限类型。通过设置粘滞位，使得只有目录的所有者才能够删除目录中的文件，即使用户组和其他用户拥有足够的权限也不能执行该删除操作。

要设置粘滞位，利用chmod将+t应用于目录：

```
$ chmod a+t directory_name
```

3. 以递归的方式设置权限

有时候，要以递归的方式修改当前目录下的所有文件和子目录的权限，可以使用下面的方法：

```
$ chmod 777 . -R
```

选项-R指定以递归的方式修改权限。

我们用“.”指定当前工作目录，这等同于：

```
$ chmod 777 "$(pwd)" -R.  
Sarath Lakshman 7 January 2011 8:41 PM
```

4. 以递归的方式设置所有权

用chown命令结合-R就可以以递归的方式设置所有权：

```
$ chown user.group . -R
```

5. 以不同的用户运行可执行文件

一些可执行文件需要以不同的用户身份（除启动该文件的当前用户之外的用户），用文件路径来执行（如 /executable_name）。有一个叫做setuid的特殊文件权限，它允许其他用户以文件所有者的身份来执行文件。

首先将该文件的所有权替换为该用户，这项操作每次都会执行，使该用户能以文件所有者的身份登录。然后运行下面的命令：

```
$ chmod +s executable_file  
  
# chown root.root executable_file  
# chmod +s executable_file  
$ ./executable_file
```

现在，这个文件实际上每次都是以超级用户的身份来执行。

setuid的使用不是无限制的。为了确保安全，它只能应用在Linux ELF格式二进制文件上，而不能用于脚本文件。

3.7 创建不可修改文件

在常见的Linux扩展文件系统中（如ext2、ext3、ext4等），可以将文件设置为不可修改（immutable）。某些文件属性可帮助我们将文件设置不可修改。一旦文件被设置为不可修改，任何用户包括超级用户都不能删除该文件，除非其不可修改的属性被移除。通过查看/etc/mstab文件，很容易找出所有挂载分区的文件系统类型。这个文件的第一列指定了分区设备路径（如/dev/sda5），第三列指定了文件系统类型（如ext3）。让我们来看看如何将文件设置为不可修改。

3.7.1 预备知识

可以用chattr将文件设置为不可修改。不过chattr能够更改的扩展属性可不止这些。

不可修改属性是保护文件不被修改的安全手段之一。最有代表性的例子就是/etc/shadow文件。该文件由当前系统中所有用户的加密密码组成。我们通过密码才能够登录系统。用户通常用passwd命令修改自己的密码。执行passwd时，它实际上就修改了/etc/shadow文件。我们可以将shadow文件设置为不可修改，这样就再没有用户能够修改密码了。让我们来看看这是如何实现的。

3.7.2 实战演练

可以按照下面的方式将一个文件设置为不可修改：

```
chattr +i file
```

或者

```
$ sudo chattr +i file
```

这样文件file就变为了不可修改。来试试下面的命令：

```
rm file  
rm: cannot remove 'file': Operation not permitted
```

如果需要使文件重新可写，可以移除不可修改属性：

```
chattr -i file
```

3.8 批量生成空白文件

有时候我们可能需要生成测试样本。比如我们可能要用程序对 1000 个文件进行操作，那应怎样生成这些测试文件呢？

3.8.1 预备知识

`touch` 命令可以用来生成空白文件；如果文件存在，则可以用它修改文件的时间戳。让我们来看看如何使用这个命令。

3.8.2 实战演练

用下面的命令可以创建一个名为filename的空白文件：

```
$ touch filename
```

批量生成名字不同的空白文件：

```
for name in {1..100}.txt
do
touch $name
done
```

在上面的代码中，`{1..100}`会扩展成一个字符串“1,2,3,4,5,6,7,...100”。除了`{1..100}.txt`，我们还可以使用其他简写样式，比如 `test{1..200}.c`、`test{a..z}.txt`等。

如果文件已经存在，那么`touch`命令将会与该文件相关的所有时间戳更改为当前时间。如果我们只想更改某些时间戳，则可以使用下面的选项。

□ `touch -a` 只更改文件访问时间（access time）。

□ `touch -m` 只更改文件内容修改时间（modification time）。

除了将时间戳更改为当前时间，我们还能够为时间戳指定特定的时间和日期：

```
$ touch -d "Fri Jun 25 20:50:14 IST 1999" filename
```

`-d`使用的日期串不一定总是以同样的格式呈现。`-d`可以接受任何的标准日期格式。我们可以忽略具体时间，使用“Jan 20 2010”这种方便的日期格式。

3.9 查找符号链接及其指向目标

在类UNIX系统中，符号链接很常见，我们会碰到各种与符号链接相关的处理工作。这则攻略可能没有什么实用性，不过它给出了处理符号链接的实践方法，这些方法没准能在编写其他shell脚本时派上用场。

3.9.1 预备知识

符号链接只不过是指向其他文件的指针。它在功能上类似于Mac OS中的别名或Windows中的快捷方式。删除符号链接不会影响到原始文件。

3.9.2 实战演练

我们可以按照下面的方法创建符号链接：

```
$ ln -s target symbolic_link_name
```

例如：

```
$ ln -l -s /var/www/ ~/web
```

这个命令在已登录用户的home目录中创建了一个名为web的符号链接。这个链接指向/var/www。这些信息可以从下面的命令输出中看到：

```
$ ls web
lrwxrwxrwx 1 slynuX slynuX 8 2010-06-25 21:34 web -> /var/www
```

web -> /var/www表明web指向/var/www。

对于每一个符号链接而言，权限标记部分以字母“l”作为起始，表示这是一个符号链接。因此，为了打印出当前目录下的符号链接，可以使用下面的命令：

```
$ ls -l | grep "^l" | awk '{ print $8 }'
```

grep对ls -l的输出进行过滤，只显示以l起始的那些行。^是字符串起始标记。awk用来打印出第8列，也就是文件名部分。

另一种方法是用find打印符号链接，如下所示：

```
$ find . -type l -print
```

在上面的命令中，我们将find命令的type参数指定为“l”，告诉find只搜索符号链接文件。-print选项将符号链接列表打印到标准输出（stdout）。而“.”表示从当前目录开始搜索。

用下列命令可以打印出符号链接的指向目标：

```
$ ls -l web | awk '{ print $10 }'
/var/www
```

ls -l命令输出的每一行都对应一个文件的详细信息。ls -l web会将符号链接文件web的详细信息全部列出。输出的第10列包含文件指向的目标（如果这是一个符号链接的话）。因此，为了找出与符号链接相关联的目标，我们可以用awk打印出文件详细列表（ls -l的输出）的第10列。

另外，我们也可以用readlink命令来完成同样的任务。这是你最应该优先考虑的方法，其使用方式如下：

```
$ readlink web
/var/www
```

3.10 列举文件类型统计信息

提到文件类型，可谓数量繁多。如果编写一个脚本，使它能够遍历一个目录中所有的文件，并生成一份关于文件类型细节以及每种文件类型数量的报告，这肯定很有意思。这则攻略正是练习如何编写一个能够遍历大量文件并收集详细信息的脚本。

3.10.1 预备知识

find命令可以通过查看文件内容来找出特定类型的文件。在UNIX/Linux系统中，文件类型并不是由文件扩展名决定的（而在微软Windows平台中正是这么做的）。编写这个脚本的目的是从多个文件中收集文件类型统计信息。我们用关联数组存储相同类型的文件数量，用find命令获取每一个文件的类型细节。

3.10.2 实战演练

用下面的命令打印文件类型信息：

```
$ file filename
```

```
$ file /etc/passwd
/etc/passwd: ASCII text
```

打印不包括文件名在内的文件类型信息：

```
$ file -b filename
ASCII text
```

生成文件统计信息的脚本如下：

```
#!/bin/bash
# 文件名: filestat.sh

if [ $# -ne 1 ];
then
    echo $0 basepath;
    echo
fi
path=$1

declare -A statarray;

while read line;
do
    ftype=`file -b "$line"`
    let statarray["$ftype"]++;
done< <(find $path -type f -print)

echo ===== File types and counts =====
for ftype in "${!statarray[@]}";
do
    echo $ftype : ${statarray["$ftype"]}
done
```

用法如下：

```
$ ./filestat.sh /home/slynux/temp
```

一个简单的输出如下：

```
$ ./filetype.sh /home/slynux/programs
----- File types and counts -----
Vim swap file : 1
ELF 32-bit LSB executable : 6
ASCII text : 2
ASCII C program text : 10
```

3.10.3 工作原理

在脚本中声明了一个关联数组`statarray`，这样可以用文件类型作为数组索引，将每种文件类型的数量存入数组。每次遇到一个文件类型，就用`let`增加计数。`find`命令以递归的方式获取文件路径列表。脚本中的`fotype=file -b "$line"`使用`file`命令获得文件类型信息。选项`-b`告诉`file`命令只打印文件类型（不包括文件名）。输出的文件类型信息包含很多细节，比如图像编码以及分辨率（如果是图像文件的话）。不过，对于这些细节我们并不感兴趣，我们只需要基本的信息就够了。各种细节信息是由逗号分隔的，例如：

```
$ file a.out -b
ELF 32-bit LSB executable, Intel 80386, version 1 (SYSV), dynamically
linked (uses shared libs), for GNU/Linux 2.6.15, not stripped
```

我们只需要从上面这些细节中提取“ELF 32-bit LSB executable”。因此我们使用`cut -d, -f1`，指明以逗号作为定界符，并且只打印第一个字段。

`done<<(find $path -type f -print);`是一段很重要的代码。它的执行逻辑如下：

```
while read line;
do something
done< filename
```

我们不用`filename`，而是用`find`命令的输出。

`<(find $path -type f -print)`等同于文件名。只不过它用子进程输出来代替文件名。注意这里还有另外一个`<`。

`${!statarray[@]}`用于返回一个数组索引列表。

3.11 环回文件与挂载

环回（loopback）文件系统是Linux系统中非常有趣的部分。我们通常是在设备上（例如磁盘分区）创建文件系统。这些存储设备能够以设备文件的形式来使用，比如`/dev/device_name`。为了使用存储设备上的文件系统，我们需要将其挂载到一些被称为挂载点（mount point）的目录上。环回文件系统是指那些在文件中而非物理设备中创建的文件系统。我们可以将这些文件挂载到挂载点上，就像设备一样。让我们来看看这是如何实现的。

3.11.1 预备知识

环回文件系统存在于文件之中。我们通过将环回文件连接到一个设备文件来进行挂载(mount)。环回文件系统的一个例子就是初始化内存文件，它位于 /boot/initrd.img。这个文件中存储了一个用于内核的初始化文件系统。

让我们来看看如何在一个1GB的文件中创建ext4文件系统。

3.11.2 实战演练

下面的命令可以创建一个1GB大小的文件。

```
$ dd if=/dev/zero of=loopbackfile.img bs=1G count=1
1024+0 records in
1024+0 records out
1073741824 bytes (1.1 GB) copied, 37.3155 s, 28.8 MB/s
```

你会发现创建好的文件超过了1GB。这是因为硬盘作为块设备，其分配存储空间时是按照块大小的整数倍来进行的。

用mkfs命令格式化这个1GB的文件：

```
# mkfs.ext4 loopbackfile.img
```

这个命令将文件格式化成ext4文件系统。使用下面的命令检查文件系统：

```
$ sudo file loopbackfile.img

loopbackfile.img: Linux rev 1.0 ext4 filesystem data, UUID=c9d56c42-f8e6-
4cbd-a9ab-369d5056660a (extents) (large files) (huge files)
```

现在就可以挂载环回文件了：

```
$ sudo mkdir /mnt/loopback
# mount -o loop loopback.img /mnt/loopback
```

o loop用来挂载环回文件系统。

这是一种快捷的挂载方法。我们并没有连接到任何设备上。但是在内部，这个环回文件连接到了一个名为/dev/loop1或loop2的设备上。

我们可以手动来操作：

```
# losetup /dev/loop1 loopback.img
# mount /dev/loop1 /mnt/loopback
```

第一种方法并不能适用于所有情况。假如我们想创建一个硬盘文件，然后对它分区并挂载其中某个分区，那就不能使用mount -o loop，而要使用第二种方法。对一个用0填充的文件进行分区：

```
# losetup /dev/loop1 loopback.img
# fdisk /dev/loop1
```

在lookback.img中创建分区并挂载第一个分区：

```
# losetup -o 32256 /dev/loop2 loopback.img
```

现在，/dev/loop2就代表第一个分区。

-o表示偏移量。32256字节是针对DOS分区方案的一个设置^①。第一个分区自硬盘第32256字节之后起始。

我们可以指定所需的偏移量来设置第二个分区。挂载过分区之后，就可以像在物理设备上面那样执行任何操作了。

如果要卸载（umount），使用下面的方法：

```
# umount mount_point
```

例如：

```
# umount /mnt/sd1
```

或者，也可以用设备文件的路径作为umount命令的参数：

```
# umount /dev/sd1
```

因为umount是一个特权命令，所以必须以超级用户的身份来执行。

3.11.3 补充内容

让我们再来研究其他的mount选项。

1. 将ISO文件作为环回文件挂载

ISO文件是任意种类光学存储介质的归档。我们可以采用环回挂载的方法，像挂载物理光盘一样挂载ISO文件。

挂载点就是一个目录，用来作为通过文件系统访问设备内容的路径。我们甚至可以用一个非空的目录作为挂载路径。那么直到设备被卸载，这个挂载路径中包含的都是来自设备的数据，而非原始目录中的内容。例如：

```
# mkdir /mnt/iso
# mount -o loop linux.iso /mnt/iso
```

现在就可以用/mnt/iso中的文件进行操作了。ISO是一个只读文件系统。

2. 通过sync立即应用更改

当对挂载设备作出更改之后，这些改变并不会被立即写入物理设备。只有当缓冲区被写满之后才会进行设备回写。但是我们可以用sync命令强制将它即刻写入更改：

```
# sync
```

sync命令必须以超级用户身份来执行。

^① losetup 中的 -o 32256 (512*63=32256) 用于设置数据偏移。由于历史原因，硬盘第一个扇区 (512 字节) 作为 MBR (Master Boot Record, 主引导记录)，其后的62个扇区作为保留扇区。

3.12 生成 ISO 文件及混合 ISO

ISO 镜像是一种存档格式，它存储了如 CD ROM、DVD ROM 等光盘的精确存储镜像内容。我们通常都是用 ISO 镜像刻录光盘。但如果你想创建一个光盘的镜像，又该如何做呢？为此，我们需要制作一个光盘的 ISO 镜像文件。很多人都是依赖于第三方工具来创建 ISO 镜像。其实，若使用命令行的话，这也就是一行命令就可以搞定的事儿。

很多人并不知道可引导光盘与不可引导光盘之间的区别。可引导光盘自身具备引导能力，也可以运行操作系统或其他软件。不可引导光盘则做不到这些。人们通常会做的一件事是：将可引导光盘中的内容复制到另一个位置作为副本，然后用副本刻录 CD。但是这样刻出的光盘并没有引导能力。要想保留光盘的可引导性，应该将它以磁盘镜像或是 ISO 文件的形式进行复制。

现在，多数人会用闪存或硬盘作为光盘的代替品。当我们将一个可引导的 ISO 文件写入闪存后，它却再也无法引导了，除非我们使用一种专门设计用于闪存设备的混合 ISO 镜像。

这则攻略将带你认识 ISO 镜像以及它的处理方法。

3.12.1 预备知识

本书已经提到过多次，UNIX 是将一切都当做文件来处理的。一个设备也是一个文件。那应该怎样复制一个设备的精确镜像呢？那就需要读出所有的数据，并将其写入另外一个文件，对吧？

正如我们所知，cat 命令可以用来读取任何数据，重定向可以用来写入文件。

3.12.2 实战演练

用下面的命令从 /dev/cdrom 创建一个 ISO 镜像：

```
# cat /dev/cdrom > image.iso
```

这条命令完全可以奏效，它从设备中读出所有的数据并将其写入一个 ISO 镜像。

用 cat 命令创建 ISO 镜像算是个小技巧。不过，创建 ISO 镜像最好的方法是使用 dd 工具。

```
# dd if=/dev/cdrom of=image.iso
```

mkisofs 命令用于创建 ISO 文件系统。可以用 cdrecord 之类的工具将 mkisofs 的输出文件直接刻录到 CD ROM 或 DVD ROM 上。我们可以将需要的所有文件放入同一个目录中，然后用 mkisofs 将整个目录的内容写入一个 ISO 文件。方法如下：

```
$ mkisofs -V "Label" -o image.iso source_dir/
```

mkisofs 命令中的选项 -o 指定了 ISO 文件的路径。source_dir 是作为 ISO 文件内容来源的目录路径，选项 -v 指定了 ISO 文件的卷标。

3.12.3 补充内容

让我们再多学一些有关 ISO 文件的命令和技术。

1. 能够启动闪存或硬盘的混合ISO

通常情况下,不能将可引导的ISO文件写入USB设备来引导操作系统。但是有一种被称为“混合ISO”的特殊的ISO文件可以做到这一切。

我们可以用isohybrid命令把标准ISO文件转换成混合ISO。isohybrid是一个比较新的工具,大多数的Linux发行版中还未包含这个工具。你可以从<http://syslinux.zytor.com>下载syslinux软件包。

来看看下面的命令:

```
# isohybrid image.iso
```

执行这个命令,我们将获得一个名为image.iso的混合ISO,它可用于写入USB存储设备。

将ISO写入USB存储设备:

```
# dd if=image.iso of=/dev/sdb1
```

你可以用适当的设备代替sdb1。

或者使用cat命令:

```
# cat image.iso > /dev/sdb1
```

2. 用命令行刻录ISO

cdrecord命令可以用来将ISO文件刻入CD ROM或DVD ROM。刻录CD ROM的方法如下:

```
# cdrecord -v dev=/dev/cdrom image.iso
```

还有一些其他的选项,如下所示。

❑ 我们可以用-speed选项指定刻录速度:

```
-speed SPEED
```

例如:

```
# cdrecord -v dev=/dev/cdrom image.iso -speed 8
```

参数8指定其刻录速度为8x。

❑ 刻录CD ROM时也可以采用多区段(multisession)方式,这样就能在一张光盘上分多次刻录数据。多区段刻录需要使用-multi选项:

```
# cdrecord -v dev=/dev/cdrom image.iso -multi
```

3. 玩转CD ROM托盘

想找点乐子,不妨试试下面的命令吧。

❑ \$ eject

这个命令用以弹出光驱托盘。

❑ \$ eject -t

这个命令用以合上光驱托盘。

不妨试着写一个可以让托盘开关N次的循环吧。

3.13 查找文件差异并进行修补

当一个文件有多个版本时,如果能够重点标记出这些版本之间的不同而无须通过人工查看来比较,那就简直是太棒了。如果文件有1000行,光靠人工进行比较可是件极其费时费力的活儿。这则攻略为你演示如何用行号对文件之间的差异进行重点标记。当多名开发人员在文件繁多的环境下工作时,如果某个人对其中某个文件进行了修改,那么应该这个修改告知其他所有开发人员。如果发送整个源代码,不仅浪费磁盘存储空间,而且单靠手动检查这些修改就很耗时间。这时,发送一个差异文件就显得很有用了。它只包含那些被修改过、添加或删除的行以及行号。这个差异文件被称为修补文件(patch file)。我们可以用patch命令将修补文件中包含的更改信息应用到原始文件中。我们也可以再次进行修补来撤销改变。来看看如何才能实现这一切。

3

3.13.1 实战演练

diff命令可以生成差异文件。

为了生成差异信息,先创建下列文件。

□ 文件1: version1.txt

```
this is the original text
line2
line3
line4
happy hacking !
```

□ 文件2: version2.txt

```
this is the original text
line2
line4
happy hacking !
GNU is not UNIX
```

非一体化(non-unified)形式的diff输出(在不使用-u选项的情况下)如下:

```
$ diff version1.txt version2.txt
3d2
<line3
6c5
> GNU is not UNIX
```

一体化形式的diff输出如下:

```
$ diff -u version1.txt version2.txt
--- version1.txt      2010-06-27 10:26:54.384884455 +0530
+++ version2.txt      2010-06-27 10:27:28.782140889 +0530
@@ -1,5 +1,5 @@
this is the original text
line2
-line3
line4
```

```
happy hacking !
-
+GNU is not UNIX
```

选项 `-u` 用于生成一体化输出。因为一体化输出的可读性更好，更易于看出两个文件之间的差异，所以人们往往更喜欢这种输出形式。

在一体化 `diff` 输出中，以 `+` 起始的是新加入的行，以 `-` 起始的是删除的行。

修补文件可以通过将 `diff` 的输出重定向到一个文件来生成：

```
$ diff -u version1.txt version2.txt > version.patch
```

现在就可以用 `patch` 命令将修改应用于任意一个文件。当应用于 `version1.txt` 时，我们就可以得到 `version2.txt`；而当应用于 `version2.txt` 时，就可以得到 `version1.txt`。

用下列命令来进行修补：

```
$ patch -p1 version1.txt < version.patch
patching file version1.txt
```

`version1.txt` 的内容现在和 `version2.txt` 的内容一模一样。

下面的命令可以撤销做出的修改：

```
$ patch -p1 version1.txt < version.patch
patching file version1.txt
Reversed (or previously applied) patch detected! Assume -R? [n] y
#修改被撤销
```

在撤销修改时，若使用 `patch` 命令的 `-R` 选项，则不会提示用户 `y/n`。

3.13.2 补充内容

让我们再看一些 `diff` 的其他特性。

生成目录的差异信息

`diff` 命令也能够以递归的形式作用于目录。它会对目录中所有的内容生成差异输出。

使用下面的命令：

```
$ diff -Naur directory1 directory2
```

上面命令中出现的选项含义如下。

- `-N`：将所有缺失的文件视为空文件。
- `-a`：将所有文件视为文本文件。
- `-u`：生成一体化输出。
- `-r`：遍历目录下的所有文件。

3.14 head 与 tail —— 打印文件的前 10 行和后 10 行

当查看上千行的大文件时，我们可不会用 `cat` 命令把整个文件内容给打印出来。相反，我们只会查看文件的一小部分内容（例如文件的前 10 行或后 10 行）。我们有可能需要打印出文件的前 `n`

行或后n行，也有可能需要打印出除了前n行或后n行之外所有的行。

还有一种用例是打印文件的第n行至第m行。

head和tail命令可以帮助我们实现这一切。

实战演练

head命令总是读取输入文件的头部。

打印前10行：

```
$ head file
```

从stdin读取数据：

```
$ cat text | head
```

指定打印前几行：

```
$ head -n 4 file
```

这个命令会打印文件的前4行。

打印除了最后N行之外所有的行：

```
$ head -n -N file
```

注意，-N表示一个负数。

例如，用下面的代码打印除了最后5行之外的所有行：

```
$ seq 11 | head -n -5
1
2
3
4
5
6
```

而下面的命令会打印出文件的第1行至第5行：

```
$ seq 100 | head -n 5
```

将最后几行排除在打印范围之外是head的非常重要的一个用法。可是总有人去寻找其他复杂的方法来达成这个目的。

打印文件的最后10行：

```
$ tail file
```

可以用下面的代码从stdin中读取输入：

```
$ cat text | tail
```

打印最后5行：

```
$ tail -n 5 file
```

打印除了前N行之外所有的行：

```
$ tail -n +(N+1)
```



例如，打印除前5行之外的所有行， $N+1=6$ ，因此使用下列命令：

```
$ seq 100 | tail -n +6
```

这条命令将打印出第6行至第100行。

tail命令的一个重要用法是从一个内容不断增加的文件中读取数据。新增加的内容总是被添加到文件的尾部，因此当新内容被写入文件的时候，可以用tail将其显示出来。只是简单地使用tail的话，它会读取文件的最后10行，然后退出。然而那时新的内容也许又已经被其他进程追加到文件中了。为了能够不间断地监视文件的增长，tail有一个特殊的选项 `-f` 或 `--follow`，它们会使tail密切关注文件中新添加的内容，并随着数据的增加时时保持更新：

```
$ tail -f growing_file
```

这种文件内容会不断增加的典型例子就是日志文件。用于监视文件内容增加的命令如下：

```
# tail -f /var/log/messages
```

或者

```
$ dmesg | tail -f
```

我们不时会运行dmesg查看内核的环形缓冲区（kernel ring buffer）消息，要么是调试USB设备，要么是查看sdX（X是sd设备的次序列号）。tail -f也可以加入一个睡眠间隔 `-s`，这样我们就可以设置监视文件更新的时间间隔。

tail具有一个很有意思的特性：当某个给定进程结束之后，tail也会随之终结。

假设我们正在读取一个不断增长的文件，进程foo一直在向该文件追加数据，那么tail -f就会一直执行，直到进程foo结束。

```
$ PID=$(pidof foo)
$ tail -f file --pid $PID
```

当进程foo结束之后，tail也会跟着结束。

让我们实际演练一下。

创建一个新文件file.txt，用gedit打开这个文件（你也可以使用其他文本编辑器）。

向文件添加新内容并不断地将它们保存到gedit中。

现在运行：

```
$ PID=$(pidof gedit)
$ tail -f file.txt --pid $PID
```

当你不断对文件进行更新时，更新的内容都会被tail命令写入终端；当关闭gedit时，tail命令也会结束。

3.15 只列出目录的其他方法

尽管看起来只列出目录是一件挺简单的任务，很多人却未必能够办到。我经常会碰到这种情况，甚至有些编写shell脚本水平不错的用户也会被这个问题难住。这则攻略很值得学习，因为它介绍了多种只列出目录的方法与技巧。

3.15.1 预备知识

有很多种方法可以只列出目录。当你向其他人询问这些技术的时候，你得到的第一个答案可能是`dir`。但是，这是错的。`dir`命令只不过是另一个类似于`ls`且比`ls`选项更少的命令。让我们来看看如何只列出目录。

3.15.2 实战演练

有4种方法可以列出当前路径下的目录，分别列出如下。

❑ `$ ls -d */`

只有上面这种结合`-d`的用法才能够只打印出目录。

❑ `$ ls -F | grep "/$"`

当使用`-F`时，所有的输出项都会添加上一个代表文件类型的字符，如`@`、`*`、`|`等。对于目录项，添加的是`/`字符。我们用`grep`只过滤那些以`/`作为行尾标记的输出项。

❑ `$ ls -l | grep "^d"`

`ls -l`输出的每一行的第一个字符表示文件类型。目录的文件类型字符是`"d"`。因此我们用`grep`过滤以`"d"`起始的行。`^`是行首标记。

❑ `$ find . -type d -maxdepth 1 -print`

`find`命令可以指定`type`的参数为目录并将`maxdepth`设置成1，这是因为我们不需要向下搜索。

3.16 在命令行中用 pushd 和 popd 快速定位

当在终端或shell提示符下涉及多处位置时，我们通常所做的就是复制并粘贴路径。有鼠标，复制并粘贴这招才管用。当没有GUI（图形用户界面），只能通过命令行进行访问的时候，就很难处理涉及在多条路径之间跳转的操作。举例来说，假如我们同时涉及`/var/www`、`/home/slynux`和`/usr/src`，当要在这些位置之间进行切换时，每次都要通过键盘输入路径，这实在是一件很麻烦的事。此时，我们就可以使用诸如`pushd`和`popd`这种基于命令行接口（CLI）的定位技术。让我们来看看它们的使用方法。

3.16.1 预备知识

`pushd`和`popd`可以用于在多个目录之间进行切换而无需复制并粘贴目录路径。`pushd`和`popd`是以栈的方式来运作。我们都知道栈是一个后进先出（Last In First Out, LIFO）的数据结构。目录路径被存储在栈中，然后用`push`和`pop`操作在目录之间进行切换。

3.16.2 实战演练

使用`pushd`和`popd`的时候，就可以无视`cd`命令了。

为了压入并切换路径，使用：

```
~ $ pushd /var/www
```

现在，栈中包含了 /var/www ~，当前目录切换到 /var/www。

然后，再压入下一个目录路径：

```
/var/www $ pushd /usr/src
```

现在栈包含 /usr/src /var/www ~，当前目录为 /usr/src。

你可以按照上面的方法，根据需要压入任意多的目录路径。

用下面的命令查看栈内容：

```
$ dirs
/usr/src /var/www ~ /usr/share /etc
0        1        2 3        4
```

当你想切换到列表中任意一个路径时，将每条路径从0到n进行编号，然后使用你希望切换到的路径编号，例如：

```
$ pushd +3
```

这条命令会将栈进行翻转并切换到目录 /usr/share。

pushd总是将路径添加到栈，如果要从栈中删除路径，可以使用popd。

移除最近压入栈的路径并切换到下一个目录：

```
$ popd
```

假设现在栈包含 /usr/src /var/www ~ /usr/share /etc，当前目录是 /usr/share，popd会将栈更改为 /var/www ~ /usr/share /etc，并且把目录切换到/var/www。

用popd +no可以从列表中移除特定的路径。

在这里，no是从左到右，从0到n开始计数的。

3.16.3 补充内容

让我们再看一些基本的目录定位练习。

在常用的目录之间切换

当涉及3个以上的目录时，可以使用pushd和popd。但是当你只涉及两个位置的时候，还有另一个更简便的方法：cd -。

如果当前路径是 /var/www，执行下面的命令：

```
/var/www $ cd /usr/src
/usr/src $ # 做点什么
```

现在要切换回 /var/www，你不需要再输入一次，只需要执行：

```
/usr/src $ cd -
```

你还可以再切换到 /usr/src：

```
/var/www $ cd -
```

3.17 统计文件的行数、单词数和字符数

在文本处理工作中，统计文件的行数、单词数和字符数非常有用。很多时候，单词或字符计数被作为一种间接的技巧来生成所需要的输出样式及结果。本书在其他章节就包含了一些这样的实例。对开发人员来说，统计LOC（Line of Code，代码行数）是一件重要的工作。我们可能需要对除无关文件之外的特定类型的文件进行统计。将wc结合其他一些命令会有助于我们实现这些需求。

3.17.1 预备知识

wc是一个用于统计的工具。它是Word Count（单词统计）的缩写。让我们来看看如何使用wc统计文件的行数、单词数和字符数。

3.17.2 实战演练

统计行数：

```
$ wc -l file
```

如果需要将stdin作为输入，使用下列命令：

```
$ cat file | wc -l
```

统计单词数：

```
$ wc -w file
```

```
$ cat file | wc -w
```

统计字符数：

```
$ wc -c file
```

```
$ cat file | wc -c
```

举个例子，我们可以按照下面的方法统计文本中的字符数：

```
echo -n 1234 | wc -c
```

```
4
```

-n用于避免echo添加额外的换行符。

当不使用任何选项执行wc时：

```
$ wc file
```

它会打印出文件的行数、单词数和字符数，彼此之间用制表符分隔。

3.17.3 补充知识

接着看看wc命令的其他选项。

打印最长行的长度

wc可以借助-L选项打印最长行的长度：

```
$ wc file -L
```

3.18 打印目录树

当你准备教程和文档时，将目录和文件系统以图形化的树状层次结构描述，一定会为你增色不少。在编写一些监控脚本时，用清晰易懂的树状描述方式来查看文件系统，同样颇有益处。让我们来看看这是如何实现的。

3.18.1 预备知识

`tree`命令是以图形化的树状结构打印文件和目录的主角。通常，在Linux发行版中并没有包含这个命令。你需要用包管理器自行安装。

3.18.2 实战演练

下面是树状UNIX文件系统的一个示例：

```
$ tree -/unixfs
unixfs/
|-- bin
|   |-- cat
|   |-- ls
|-- etc
|   |-- passwd
|-- home
|   |-- pactpub
|   |   |-- automate.sh
|   |   |-- schedule
|   |-- slynux
|-- opt
|-- tmp
-- usr
8 directories, 5 files
```

`tree`命令有很多有意思的选项，让我们认识一下其中的几项。

只重点标记出匹配某种样式的文件：

```
$ tree path -P PATTERN #用通配符描述样式
```

例如：

```
$ tree PATH -P "*.sh" #用一个目录路径代替PATH
|-- home
|   |-- pactpub
|   |   |-- automate.sh
```

只重点标记出除符合某种样式之外的那些文件：

```
$ tree path -I PATTERN
```

为了同时打印出文件和目录的大小，可以使用 `-h` 选项：

```
$ tree -h
```



3.18.3 补充内容

请看一个有趣的tree命令选项。

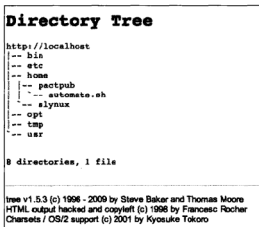
以HTML形式输出目录树

tree命令可以生成HTML输出。例如，用下面的命令可以创建一个包含目录树输出的HTML文件。

```
$ tree PATH -H http://localhost -o out.html
```

将http://localhost替换为适合存放输出文件的URL。将PATH替换为主目录（base directory）的真正路径。对当前目录可以用“.”作为PATH。

根据目录列表生成的Web页面形式如图3-1所示。



```
Directory Tree

http://localhost
|-- bin
|-- etc
|-- home
|   |-- pactpub
|   |-- automate.sh
|   |-- slynux
|-- opt
|-- tmp
|-- usr

8 directories, 1 file

tree v1.5.3 (c) 1996 - 2009 by Steve Baker and Thomas Moore
HTML output hacked and copyleft (c) 1996 by Francesc Rocher
Charsets / OS/2 support (c) 2001 by Kyosuke Tokoro
```

图 3-1

本章内容

- 正则表达式入门
- 用grep在文件中搜索文本
- 用cut按列切分文件
- 统计特定文件中的词频
- sed入门
- awk入门
- 替换文本或文件中的字符串
- 压缩或解压缩JavaScript
- 对文件中的行、单词和字符进行迭代
- 按列合并文件
- 打印文件或行中的第 n 个单词或列
- 打印不同行或样式之间的文本
- 用脚本检查回文字符串
- 以逆序形式打印行
- 解析文本中的电子邮件地址和URL
- 打印文件中某个样式之前或之后的 n 行
- 在文件中移除包含某个单词的句子
- 用awk实现head、tail和tac
- 文本切片与参数操作

4.1 简介

shell脚本语言包含了众多用于解决UNIX/Linux系统问题必不可少的组件。对于UNIX环境中出现的难题，Bash总是能够给出快捷的解决之道。文本处理是shell脚本的重要应用领域之一。shell脚本可以将sed、awk、grep、cut等这类优美的工具组合在一起，用于解决文本处理相关的问题。大多数编程语言都被设计成了通用语言，因此要编写一个能够处理文本并生成理想输出的程序，就不得不花费大量的气力。由于在设计Bash之初设计人员就考虑到了文本处理，Bash得以包含大量的便于文本处理的功能。

有各种各样的工具能够帮助我们从不同的细节层面上处理文件，比如字符、行、单词、行列等。因此我们可以用很多方法来处理一个文本文件。正则表达式是样式匹配技术的核心。借助适合的正则表达式，可以生成我们所需的如过滤、剥离（stripping）、替换、搜索等各类输出结果。

本章包括一系列攻略，探讨了多个与文本处理相关的问题，了解这些问题在编写用于解决实际问题的脚本时大有裨益。

4.2 正则表达式入门

正则表达式是基于样式匹配的文本处理技术的关键所在。想要在编写文本处理工具方面驾轻就熟，你就必须对正则表达式有一个基本的理解。正则表达式是一种用于文本匹配的形式小巧、具有高度针对性的编程语言。只依靠通配符技术，能够匹配的文本范围相当有限。这则攻略将对基础的正则表达式进行讲解。

4.2.1 预备知识

正则表达式是用于绝大多数文本处理工具的一种语言。因此你完全可以把在这里学习到的技术应用到其他攻略中。`[a-z0-9_]+@[a-z0-9]+\.[a-z]+`就是一个能够匹配电子邮件地址的正则表达式。

它看起来是不是很奇怪？别担心，一旦你理解了相关的概念，它其实相当简单。

4.2.2 实战演练

本节将会让你简单了解一下正则表达式、POSIX字符类（POSIX character class）以及元字符（meta character）。

先来看一看正则表达式的基本组成部分，如表4-1所示。

表 4-1

正则表达式	描 述	示 例
<code>^</code>	行起始标记	<code>^tux</code> 匹配以tux起始的行
<code>\$</code>	行尾标记	<code>tux\$</code> 匹配以tux结尾的行
<code>.</code>	匹配任意一个字符	<code>Hack.</code> 匹配Hack1和Hacki，但是不能匹配Hack12和Hackil，它只能匹配单个字符
<code>[]</code>	匹配包含在[字符]之中的任意一个字符	<code>coo[kl]</code> 匹配cook或cool
<code>[^]</code>	匹配除[字符]之外的任意一个字符	<code>9[^01]</code> 匹配92、93，但是不匹配91或90
<code>[-]</code>	匹配[]中指定范围内的任意一个字符	<code>[1-5]</code> 匹配从1~5的任意一个数字
<code>?</code>	匹配之前的项1次或0次	<code>colou?r</code> 匹配color或colour，但是不能匹配colour
<code>+</code>	匹配之前的项1次或多次	<code>Rollno-9+</code> 匹配Rollno-99、Rollno-9，但是不能匹配Rollno-
<code>*</code>	匹配之前的项0次或多次	<code>co*l</code> 匹配cl、col、cool等
<code>()</code>	创建一个用于匹配的子串	<code>ma(tri)?</code> 匹配max或matrix
<code>{n}</code>	匹配之前的项n次	<code>[0-9]{3}</code> 匹配任意一个三位数， <code>[0-9]{3}</code> 可以扩展为 <code>[0-9][0-9][0-9]</code>
<code>{n,}</code>	之前的项至少需要匹配n次	<code>[0-9]{2,}</code> 匹配任意一个两位或更多位的数字
<code>{n,m}</code>	指定之前的项所必需匹配的最小次数和最大次数	<code>[0-9]{2,5}</code> 匹配从两位数到五位数之间的任意一个数字
<code> </code>	交替——匹配 两边的任意一项	<code>Oct (1st 2nd)</code> 匹配Oct 1st或Oct 2nd
<code>\</code>	转义符可以将上面介绍的特殊字符进行转义	<code>a\.b</code> 匹配ab，但不能匹配ajb。通过在. 之间加上前缀\，从而忽略了. 的特殊意义

POSIX字符类是一个形如[:...:]的特殊元序列 (meta sequence)，它可以用于匹配特定的字符范围。POSIX字符类如表4-2所示。

表 4-2

正则表达式	描 述	示 例
[[:alnum:]]	字母与数字字符	[[:alnum:]]+
[[:alpha:]]	字母字符 (包括大写字母与小写字母)	[[:alpha:]]{4}
[[:blank:]]	空格与制表符	[[:blank:]]*
[[:digit:]]	数字字符	[[:digit:]]?
[[:lower:]]	小写字母	[[:lower:]]{5,}
[[:upper:]]	大写字母	([[:upper:]]+)?
[[:punct:]]	标点符号	[[:punct:]]
[[:space:]]	包括换行符、回车等在内的所有空白字符	[[:space:]]+

元字符是一种Perl风格的正则表达式，只有一部分文本处理工具支持它，并不是所有的工具都支持表4-3中所列的字符，但是之前介绍的正则表达式和字符类都是被广泛支持的。

表 4-3

正则表达式	描 述	示 例
\b	单词边界	\bcool\b匹配cool，但不匹配coolant
\B	非单词边界	cool\b匹配coolant，但不匹配cool
\d	单个数字字符	b\d匹配b2b，但不匹配bcb
\D	单个非数字字符	b\D匹配bcb，但不匹配b2b
\w	单个单词字符 (字母、数字与 _)	\w匹配l或a，但不匹配&
\W	单个非单词字符	\W匹配&，但不匹配l或a
\n	换行符	\n匹配一个换行
\s	单个空白字符	x\s匹配x x，但不匹配xx
\S	单个非空白字符	\x\S匹配xxx，但不匹配xx
\r	回车	\r匹配回车

4.2.3 工作原理

在4.2.2节中看到的表格是正则表达式的关键元素表。在表中挑选恰当的正则表达式，我们就能够根据需要构建出适合的正则表达式来匹配文本。正则表达式是一种用于文本匹配的通用语言。因此我们在这则攻略中不会介绍任何工具，而是把这个任务放到本章的其他攻略中。

来看几个文本匹配的例子。

□ 为了匹配给定文本中的所有单词，可以使用下面的正则表达式：

```
(?[a-zA-Z])+ ?)
```

“?”用于匹配单词前后可能出现的空格。[a-zA-Z]+代表一个或多个字母 (a-z和A-Z)。

□ 为了匹配一个IP地址，可以使用下面的正则表达式：

```
[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}\.[0-9]{1,3}
```

或者

```
[[[:digit:]]{1,3}\.[:digit:]]{1,3}\.[:digit:]]{1,3}\.[:digit:]]{1,3}
```

我们知道IP地址通常的书写形式是192.168.0.2，它是由点号分割的4个整数（每一个整数的取值范围从0到255）。

[0-9]或[:digit:]匹配数字0~9。{1,3}匹配1到3个数字，\匹配"。"。

4.2.4 补充内容

下面就看一看如何在正则表达式中指定某些字符的特殊含义。

处理特殊字符

正则表达式用\$、^、.、*、+、{以及}等作为特殊字符。但是如果我们希望将这些字符作为非特殊字符（表示普通字面含义的字符）来使用的话，应该怎么做呢？来看一个例子。

正则表达式：`[a-z]*.[0-9]`

那么，它是什么意思？

它可以匹配0个或多个[\[a-z\]](#)（`[a-z]*`），接任意单个字符（`.`），再接[\[0-9\]](#)中的任意一个字符，所以它能够匹配 `abcde09`。

它也可以理解成：匹配[\[a-z\]](#)中任意一个字符，接单个字符`*`，再接单个字符（点号），最后接一个数字，所以它能够匹配`x*.8`。

为了避免这种理解上的混乱，我们可以在字符前面放上一个“\”（这种做法称为“对字符进行转义”）。对于像`*`这种具有多种含义的字符，可以在前面加上“\”，使其具备或丧失某些特殊的含义。至于转义后字符的意义是否具备特殊的含义，则取决于你所使用的工具。

4.3 用 grep 在文件中搜索文本

在文件中进行搜索是文本处理中一项重要的用例。我们也许需要根据某些规格在上千行的文件中查找所需要的数据。这则攻略将教你如何从大量数据中定位符合给定规格的数据。

4.3.1 预备知识

`grep`命令是UNIX中用于文本搜索的大师级工具。它能够接受正则表达式和通配符。我们可以用多个`grep`命令选项来生成各种格式的输出。让我们看看它具体的使用方法。

4.3.2 实战演练

在文件中搜索一个单词：

```
$ grep match_pattern filename
this is the line containing match_pattern
```

或者

```
$ grep "match_pattern" filename
this is the line containing match_pattern
```

命令会返回包含给定match_pattern的文本行。

也可以像下面这样从stdin中读取：

```
$ echo -e "this is a word\nnext line" | grep word
this is a word
```

一个grep命令也可以对多个文件进行搜索：

```
$ grep "match_text" file1 file2 file3 ...
```

用--color选项可以在输出行中重点标记出匹配到的单词：

```
$ grep word filename --color=auto
this is the line containing word
```

grep命令通常将match_pattern视为通配符。如果要使用正则表达式，需要添加-E选项——这意味着使用扩展（extended）正则表达式，也可以使用默认允许正则表达式的grep命令——egrep。例如：

```
$ grep -E "[a-z]+"
```

或者

```
$ egrep "[a-z]+"
```

为了只输出文件中匹配到的文本部分，可以使用选项 -o：

```
$ echo this is a line. | grep -o -E "[a-z]+\."
line
```

或者

```
$ echo this is a line. | egrep -o "[a-z]+\."
line.
```

要想打印除包含match_pattern的行之外的所有行，可使用：

```
$ grep -v match_pattern file
```

选项-v可以将匹配结果进行反转（invert）。

统计文件或文本中包含匹配字符串的行数：

```
$ grep -c "text" filename
10
```

需要注意的是-c只是统计匹配行的数量，并不是匹配的次数。例如：

```
$ echo -e "1 2 3 4\nhello\n5 6" | egrep -c "[0-9]"
2
```

尽管有6个匹配项，但命令只打印出2，这是因为只有两个匹配行。在单行中出现的多次匹配

只统计为一次。

为了文件中统计匹配项的数量，可以使用下面的技巧：

```
$ echo -e "1 2 3 4\nhello\n5 6" | egrep -o "[0-9]" | wc -l
6
```

打印出包含匹配字符串的行数：

```
$ cat sample1.txt
gnu is not unix
linux is fun
bash is art
$ cat sample2.txt

planetlinux

$ grep linux -n sample1.txt
2:linux is fun
```

或者

```
$ cat sample1.txt | grep linux -n
```

如果使用了多个文件，它也会随着输出结果打印出文件名：

```
$ grep linux -n sample1.txt sample2.txt
sample1.txt:2:linux is fun
sample2.txt:2:planetlinux
```

打印样式匹配所位于的字符或字节偏移：

```
$ echo gnu is not unix | grep -b -o "not"
7:not
```

一行中字符串的字符偏移是从该行的第一个字符开始计算，起始值是0。在上面的例子中，“not”的偏移值是7[也就是说，not是从该行（即“gnu is not unix”这一行）的第7个字符开始的]。

选项 -b总是和 -o配合使用。

搜索多个文件并找出匹配文本位于哪一个文件中：

```
$ grep -l linux sample1.txt sample2.txt
sample1.txt
sample2.txt
```

和 -l相反的选项是 -L，它会返回一个不匹配的文件列表。

4.3.3 补充内容

我们已经看过grep命令的一些基本用法。不过grep的本事可不止如此，它具有各种丰富的特性。让我们接着看看grep的其他选项。

1. 递归搜索文件

如果需要在多级目录中对文本进行递归搜索，可以使用：

```
$ grep "text" . -R -n
```

命令中的“.”指定了当前目录。

例如：

```
$ cd src_dir
$ grep "test_function()" . -R -n
./miscutils/test.c:16:test_function();
```

test_function() 位于miscutils/test.c的第16行。



这是开发人员使用最多的命令之一。它用于查找某些文本位于哪些源码文件中。

2. 忽略样式中的大小写

选项 -i 可以使匹配样式不考虑字符的大小写，例如：

```
$ echo hello world | grep -i "HELLO"
hello
```

3. 用grep匹配多个样式

在进行匹配的时候通常只指定一个样式。然而，我们可以用选项 -e 来指定多个匹配样式：

```
$ grep -e "pattern1" -e "pattern"
```

例如：

```
$ echo this is a line of text | grep -e "this" -e "line" -o
this
line
```

还有另一种方法也可以指定多个样式。我们可以提供一个样式文件用于读取样式。在样式文件中逐行写下需要匹配的样式，然后用选项 -f 执行grep：

```
$ grep -f pattern_file source_filename
```

例如：

```
$ cat pat_file
hello
cool

$ echo hello this is cool | grep -f pat_file
hello this is cool
```

4. 在grep搜索中包括或排除文件

grep 可以在搜索过程中包括或排除某些文件。我们可以使用通配符来指定所需要包括或排除的文件。

只在目录中递归搜索所有的 .c 和 .cpp 文件：

```
$ grep "main()" . -r --include *.c,*.cpp
```

注意，some(string1,string2,string3) 会扩展成 somestring1 somestring2

somestring3。

在搜索中排除所有的README文件：

```
$ grep "main()" . -r --exclude "README"
```

如果需要排除目录，可以使用 `--exclude-dir` 选项。

如果需要从文件中读取所需排除的文件列表，使用 `--exclude-from FILE`。

5. 使用0值字节后缀的grep与xargs

`xargs`命令通常用于将文件名列表作为命令行参数提供给其他命令。当文件名用作命令行参数时，建议用0值字节作为文件名终止符，而不是用空格。因为一些文件名中会包含空格字符，一旦它被误解为终结符，那么单个文件名就会被认为是两个文件名（例如，New file.txt被解析成New和file.txt两个文件名）。而这个问题完全可以利用0值字节后缀来避免。我们使用`xargs`以便于从诸如`grep`、`find`等命令中接受`stdin`文本。这些命令可以将带有0值字节后缀的文本输出到`stdout`。为了指明输入的文件名是以0值字节（\0）作为终止符，我们应该在`xargs`中使用`-0`。

创建测试文件：

```
$ echo "test" > file1
$ echo "cool" > file2
$ echo "test" > file3
```

在下面的命令序列中，`grep`输出以0值字节作为终结符的文件名（\0）。这可以用`grep`的`-Z`选项来指定。`xargs -0`读取输入并用0值字节终结符分隔文件名：

```
$ grep "test" file* -lZ | xargs -0 rm
```

`-Z`通常和`-l`结合使用。

6. grep的静默输出

之前提到`grep`能够返回不同格式的输出。在某些情况下，我们需要知道一个文件是否包含指定的文本。对此，我们要执行一个可以返回真假的条件测试。这可以通过设置`grep`的静默条件（`-q`）来实现。在静默模式（quiet mode）中，`grep`命令不会向标准输出打印任何输出。它仅是运行命令，然后根据命令执行成功与否返回退出状态。

我们都知道如果命令运行成功会返回0，如果失败则返回非0值。

让我们来看一个脚本，这个脚本利用静默模式的`grep`来测试文本匹配是否存在于某个文件中。

```
#!/bin/bash
#文件名: silent_grep.sh
#用途: 测试文件是否包含特定的文本内容

if [ $# -ne 2 ];
then
echo "$0 match_text filename"
fi

match_text=$1
filename=$2

grep -q $match_text $filename
```

```
if [ $? -eq 0 ];  
then  
echo "The text exists in the file"  
else  
echo "Text does not exist in the file"  
fi
```

这个silent_grep.sh脚本可以提供用于匹配的单词(Student)和一个文件名(student_data.txt)作为命令参数:

```
$ ./silent_grep.sh Student student_data.txt  
The text exists in the file
```

7. 打印出匹配文本之前或之后的行

基于上下文的打印是grep的特色之一。假设已经找到了给定文本的匹配行,通常情况下grep只会打印出这一行。但我们也许需要匹配行之前或之后的n行,也可能两者皆要。这可以在grep中用前后行控制来实现。让我们看看具体的做法。

要打印匹配某个结果之后的3行,使用-A选项:

```
$ seq 10 | grep 5 -A 3  
5  
6  
7  
8
```

要打印匹配某个结果之前的3行,使用-B选项:

```
$ seq 10 | grep 5 -B 3  
2  
3  
4  
5
```

要打印匹配某个结果之前以及之后的3行,使用-C选项:

```
$ seq 10 | grep 5 -C 3  
2  
3  
4  
5  
6  
7  
8
```

如果有多个匹配,那么以一行"--"作为各匹配之间的定界符:

```
$ echo -e "a\nb\nc\na\nb\nc" | grep a -A 1  
a  
b  
--  
a  
b
```



4.4 用 cut 按列切分文件

我们也许不需要按行，而是需要按列切分文件。假设我们有一个文本文件，其中按行包含了学生的报表信息，例如No（学号）、Name（姓名）、Mark（成绩）和Percentage（百分比）。我们需要将学生的姓名提取到另一个文件，或者是需要表中的第 n 列或其中两列甚至更多列。这则攻略将为你演示如何完成这项任务。

4.4.1 预备知识

cut是一个帮我们将文本按列进行切分的小工具。它也可以指定分隔每列的定界符。在cut的术语中，每列都是一个字段。

4.4.2 实战演练

为了提取第一个字段或列，可以使用下面的语法：

```
cut -f FIELD_LIST filename
```

FIELD_LIST是需要显示的列。它由列号组成，彼此之间用逗号分隔。例如：

```
$ cut -f 2,3 filename
```

这条命令将显示第2列和第3列。

cut也能够从stdin中读取输入文本。

制表符是字段或列的默认定界符。没有定界符的行也会照原样被打印出来。要避免打印出这种不包含定界符的行，则可以使用cut的-s选项。一个cut命令的例子如下：

```
$ cat student_data.txt
No Name Mark Percent
1 Sarath 45 90
2 Alex 49 98
3 Anu 45 90
```

```
$ cut -f1 student_data.txt
No
1
2
3
```

提取多个字段：

```
$ cut -f2,4 student_data.txt
Name Percent
Sarath 90
Alex 98
Anu 90
```

要打印多列，需要提供一个由逗号分隔的列号列表作为-f选项的参数。

我们也可以使用-complement选项对提取的字段进行补集运算。假设有多个字段，你希望打

印出除第3列之外所有的列，则可以使用：

```
$ cut -f3 --complement student_data.txt
No Name Percent
1 Sarath 90
2 Alex 98
3 Anu 90
```

要指定字段的定界符，使用 -d 选项：

```
$ cat delimited_data.txt
No;Name;Mark;Percent
1;Sarath;45;90
2;Alex;49;98
3;Anu;45;90

$ cut -f2 -d";" delimited_data.txt
Name
Sarath
Alex
Anu
```

4.4.3 补充内容

cut 命令有一些选项可以将一串字符作为列来显示。让我们来看看这些选项。

指定字段的字符或字节范围

假设我们不依赖定界符，但需要通过将字段定义为一个字符范围（行首记为0）来进行字段提取，这种需求也可以用 cut 来实现。

表4-4中列出了字符字段范围的记法。

表 4-4

记 法	范 围
N-	从第N个字节，字符或字段到行尾
N-M	从第N个字节，字符或字段到第M个（包括第M个在内）字节，字符或字段
-M	第1个字节，字符或字段到第M个（包括第M个在内）字节，字符或字段

用上面介绍的记法，再结合下列选项将某个范围的字节或字符指定为字段：

- ☐ -b 表示字节；
- ☐ -c 表示字符；
- ☐ -f 表示定义字段。

例如：

```
$ cat range_fields.txt
abcdefghijklmnopqrstuvwxyz
abcdefghijklmnopqrstuvwxyz
abcdefghijklmnopqrstuvwxyz
abcdefghijklmnopqrstuvwxyz
```

你可以打印第1个到第5个字符；

```
$ cut -c1-5 range_fields.txt
abcde
abcde
abcde
abcde
```

打印前2个字符:

```
$ cut range_fields.txt -c-2
ab
ab
ab
ab
```

要将 `-c` 替换成 `-b`，可以用字节作为计数单位。

在使用 `-c`、`-f` 和 `-b` 时，我们可以指定输出定界符：

```
--output-delimiter 'delimiter string'
```

当用 `-b` 或 `-c` 提取多个字段时，必须使用 `--output-delimiter`，否则，你就没法区分不同的字段了。例如：

```
$ cut range_fields.txt -c1-3,6-9 --output-delimiter ","
abc,fghi
abc,fghi
abc,fghi
abc,fghi
```

4

4.5 统计特定文件中的词频

查找文件中使用单词的频率是一个很有意思的练习，在这个练习中会应用到你已学到的文本处理技巧。统计词频的方法有很多种。让我们看看具体的做法。

4.5.1 预备知识

我们可以使用关联数组、`awk`、`sed`、`grep` 等不同的方式来解决这个问题。

4.5.2 实战演练

单词是由空格和点号分隔的字母组合。首先我们应该解析出给定文件中出现的所有单词，这样才能开始统计每个单词的个数。单词解析可以用正则表达式配合 `sed`、`awk` 或 `grep` 等工具来完成。

要找出每个单词出现的次数，我们要采用一种与众不同的方法。有一种统计词频的方法是用一个循环来遍历每一个单词，然后用另一个循环来检查单词是否相同。如果相同，则增加计数并在文件尾部进行打印。这不是一种高效的方法。在关联数组中，我们将单词作为数组索引，单词计数作为数组值。通过遍历每一个单词，我们只需要使用一个循环就能完成统计任务。将数组值初始化为0，并使用 `array[word]=array[word] + 1` 就可以得到一个包含各单词出现次数的数组。

现在就来动手试试吧。创建脚本如下：

```
#!/bin/bash
#文件名: word_freq.sh
#用途: 计算文件中单词的词频

if [ $# -ne 1 ];
then
echo "Usage: $0 filename";
exit -1
fi

filename=$1

egrep -o "\b[[:alpha:]]+\b" $filename | \

awk '{ count[$0]++ }
END{ printf("%-14s%s\n","Word","Count") ;
for(ind in count)
{ printf("%-14s%d\n",ind,count[ind]); }
}'
```

输出如下:

```
$ ./word_freq.sh words.txt
Word      Count
used      1
this      2
counting  1
```

4.5.3 工作原理

`egrep -o "\b[[:alpha:]]+\b" $filename`用来只输出单词。用 `-o` 选项打印由换行符分隔的匹配字符序列。这样我们就可以在每行中列出一个单词。

`\b` 是单词边界标记符。`[[:alpha:]]` 是表示字母的字符类。

`awk` 命令用来避免对每一个单词进行迭代。因为 `awk` 默认对每一行都执行 `{}` 块中的语句, 所以我们就不需要再为同样的事手动指定循环了。借助关联数组, 当执行 `count[$0]++` 时, 单词计数就增加。最后, 在 `END{}` 语句块中通过迭代所有的单词, 就可以打印出单词及它们各自出现的次数。

4.5.4 参考

- 1.6节讲解了 Bash 中的数组。
- 4.7节介绍了 `awk` 命令。

4.6 sed 入门

`sed` 是 stream editor (流编辑器) 的缩写。它是文本处理中非常重要的工具。它能够完美地配合正则表达式使用, 功能不同凡响。`sed` 命令众所周知的一个用法是进行文本替换。这则攻略包括了 `sed` 命令大部分的常用技术。

4.6.1 实战演练

sed可以替换给定文本中的字符串。它可以利用正则表达式进行匹配。

```
$ sed 's/pattern/replace_string/' file
```

或者

```
$ cat file | sed 's/pattern/replace_string/' file
```

这个命令从stdin中读取输入。

使用-i选项,可以将替换结果应用于原文件。很多用户在进行替换之后,借助重定向来保存文件:

```
$ sed 's/text/replace/' file > newfile
$ mv newfile file
```

其实只需要一行命令就可以搞定,例如:

```
$ sed -i 's/text/replace/' file
```

之前看到的sed命令会将每一行中第一处符合样式的内容替换掉。但是如果要替换所有内容,我们需要在命令尾部加上参数g,其方法如下:

```
$ sed 's/pattern/replace_string/g' file
```

后缀/g意味着sed会替换每一处匹配。但是有时候我们不需要替换前N处匹配,而是需要替换剩下的匹配。有一个选项可以用来忽略前N处匹配,并从第N+1处开始替换。

请看下面的命令:

```
$ echo this thisthisthis | sed 's/this/THIS/2g'
thisTHISTHISTHIS
```

```
$ echo this thisthisthis | sed 's/this/THIS/3g'
thisthisTHISTHIS
```

```
$ echo this thisthisthis | sed 's/this/THIS/4g'
thisthisthisTHIS
```

当需要从第N处匹配开始替换时,可以使用 /Ng。

字符/在sed中作为定界符使用。我们可以像下面一样使用任意的定界符:

```
sed 's:text:replace:g'
sed 's|text|replace|g'
```

当定界符出现在样式内部时,我们必须用前缀\对它进行转义:

```
sed 's|site\\|xt|replace|g'
```

\\是一个出现在样式内部并经过转义的定界符。

4.6.2 补充内容

sed命令包含多个可用于文本处理的选项。将这些选项以合理的次序组合,可以在一行命令中解决很多复杂的问题。让我们看看这些选项。

1. 移除空白行

用sed移除空白行不过是小菜一碟。空白行可以用正则表达式`^$`进行匹配：

```
$ sed '/^$/d' file
```

/pattern/d会移除匹配样式的行。

在空白行中，行尾标记紧随着行首标记。

2. 已匹配字符串标记&

在sed中，用&标记匹配样式的字符串，就能够在替换字符串时使用已匹配的内容。

例如：

```
$ echo this is an example | sed 's/\w+/[&]/g'
[this] [is] [an] [example]
```

正则表达式`\w+`匹配每一个单词，然后我们用[&]替换它。&对应于之前所匹配到的单词。

3. 子串匹配标记\1

&代表匹配给定样式的字符串。但我们也可以匹配给定样式的其中一部分。来看看具体的做法。

```
$ echo this is digit 7 in a number | sed 's/digit \([0-9]\)/\1/'
this is 7 in a number
```

这条命令将digit 7替换为7。样式中匹配到的子串是7。\`(pattern\)`用于匹配子串。模式被包括在使用斜线转义过的()`()`中。对于匹配到的第一个子串，其对应的标记是`\1`，匹配到的第二个子串是`\2`，往后依次类推。下面的示例中包含了多个匹配：

```
$ echo seven EIGHT | sed 's/\([a-z]\+\)\ \([A-Z]\+\)/\2 \1/'
EIGHT seven
```

`([a-z]\+)`匹配第一个单词，`([A-Z]\+)`匹配第二个单词。`\1`和`\2`用来引用它们。这种引用被称为向后引用(back referencing)。在替换部分，它们的次序被更改为`\2 \1`，因此结果就呈现出逆序的形式。

4. 组合多个表达式

利用管道组多个sed命令这种方法可以用下面的方式代替：

```
sed 'expression' | sed 'expression'
```

它等价于：

```
$ sed 'expression; expression'
```

5. 引用

sed表达式通常用单引号来引用。不过也可以使用双引号。双引号会通过表达式求值来对其进行扩展。当我们想在sed表达式中使用一些变量字符串时，双引号就有用武之地了。

例如：

```
$ text=hello
$ echo hello world | sed "s/$text/HELLO/"
HELLO world
```

\$text的求值结果是hello。

4.7 awk 入门

作为一款工具, awk被设计用于数据流。它非常有趣, 其原因就在于它可以对列和行进行操作。awk有很多内建的功能, 比如数组、函数等, 这是它和C语言的相同之处。灵活性是awk最大的优势。

4.7.1 实战演练

awk脚本的结构基本如下所示:

```
awk ' BEGIN{ print "start" } pattern { commands } END{ print "end" }
file
```

awk命令也能够读取stdin中的内容。

一个awk脚本通常由3部分组成: BEGIN语句块、END语句块和能够使用模式匹配的通用语句块。这3个部分是可选的, 它们中任何一个部分都可以不出现在脚本中。脚本通常会被包含在单引号或双引号中:

```
awk 'BEGIN { statements } { statements } END { end statements }'
```

或者也可以使用

```
awk "BEGIN { statements } { statements } END { end statements }"
```

例如:

```
$ awk 'BEGIN { i=0 } { i++ } END{ print i}' filename
```

或者

```
$ awk "BEGIN { i=0 } { i++ } END{ print i }" filename
```

4.7.2 工作原理

awk命令的工作方式如下所示。

(1) 执行BEGIN { commands } 语句块中的语句。

(2) 从文件或stdin中读取一行, 然后执行pattern { commands }。重复这个过程, 直到文件全部被读取完毕。

(3) 当读至输入流 (input stream) 末尾时, 执行END { commands } 语句块。

BEGIN语句块在awk开始从输入流中读取行之前被执行。这是一个可选的语句块, 诸如变量初始化、打印输出表格的表头等语句通常都可以写入BEGIN语句块中。

END语句块和BEGIN语句块类似。END语句块在awk从输入流中读取完所有的行之后即被执行。像打印所有行的分析结果这类汇总信息, 都是在END语句块中实现的常见任务 (例如, 在比较过所有的行之后, 打印出最大数)。它也是一个可选的语句块。

最重要的部分就是pattern语句块中的通用命令。这个语句块同样是可选的。如果不提供该

语句块，则默认执行{ print }，即打印每一个读取到的行。awk对于读取的每一行，都会执行这个语句块。

这就像一个用来读取行的while循环，在循环体中提供了相应的语句。

每读取一行时，它就会检查该行和提供的样式是否匹配。样式本身可以是正则表达式、条件以及行匹配范围等。如果当前行匹配该样式，则执行{ }中的语句。

样式是可选的。如果没有提供样式，那么它就会默认所有的行都是匹配的，并执行{ }中的语句。

让我们看看下面的例子：

```
$ echo -e "line1\nline2" | awk 'BEGIN{ print "Start" } { print } END{
print "End" }'
Start
line1
line2
End
```

当使用不带参数的print时，它会打印出当前行。关于print，需要记住两件重要的事情：当print的参数是以逗号进行分隔时，参数打印时则以空格作为定界符；在awk的print语句中，双引号是被当做拼接操作符（concatenation operator）使用的。

例如：

```
$ echo | awk '{ var1="v1"; var2="v2"; var3="v3"; \
print var1,var2,var3 ; }'
```

该语句将按照下面的格式打印变量值：

```
v1 v2 v3
```

echo命令向标准输出写入一行，因此awk的{ }语句块中的语句只被执行一次。如果awk的标准输入包含多行，那么{ }语句块中的命令就会被执行多次。

拼接的使用方法如下：

```
$ echo | awk '{ var1="v1"; var2="v2"; var3="v3"; \
print var1"--var2"--var3 ; }'
```

输出为：

```
v1-v2-v3
```

{ }类似于一个循环体，会对文件中的每一行进行迭代。



我们通常将变量初始化语句（如var=0）以及打印文件头部的语句放入BEGIN语句块中。在END{ }语句块中，我们往往会放入打印结果等语句。

4.7.3 补充内容

awk命令具有丰富的特性。要想洞悉awk编程的精妙之处，首先应该熟悉awk重要的选项和功能。让我们来看看awk的一些重要功能。

1. 特殊变量

以下是可以用于awk的一些特殊变量。

- NR: 表示记录数量 (number of records), 在执行过程中对应于当前行号。
- NF: 表示字段数量 (number of fields), 在执行过程中对应于当前行的字段数。
- \$0: 这个变量包含执行过程中当前行的文本内容。
- \$1: 这个变量包含第一个字段的文本内容。
- \$2: 这个变量包含第二个字段的文本内容。

例如:

```
$ echo -e "line1 f2 f3\nline2 f4 f5\nline3 f6 f7" | \
awk '{
print "Line no:"NR",No of fields:"NF, "$0="$0, "$1="$1,$2="$2,$3="$3
}'
Line no:1,No of fields:3 $0=line1 f2 f3 $1=line1 $2=f2 $3=f3
Line no:2,No of fields:3 $0=line2 f4 f5 $1=line2 $2=f4 $3=f5
Line no:3,No of fields:3 $0=line3 f6 f7 $1=line3 $2=f6 $3=f7
```

我们可以用print \$NF打印一行中最后一个字段, 用 \$(NF-1)打印倒数第二个字段, 其他字段依次类推即可。

awk的printf()函数的语法和C语言中的同名函数一样。我们也可以用这个函数来代替print。

再来看awk的一些基本用法。

打印每一行的第2和第3个字段:

```
$awk '{ print $3,$2 }' file
```

要统计文件中的行数, 使用下面的命令:

```
$ awk 'END{ print NR }' file
```

这里只使用了END语句块。每读入一行, awk会将NR更新为对应的行号。当到达最后一行时, NR中的值就是最后一行的行号, 于是, 位于END语句块中的NR就包含了文件的行数。

你可以将每一行中第一个字段的值按照下面的方法进行累加:

```
$ seq 5 | awk 'BEGIN{ sum=0; print "Summation:" }
{ print $1"+", sum+=$1 } END { print "=="; print sum }'
Summation:
1+
2+
3+
4+
5+
==
15
```

2. 将外部变量值传递给awk

借助选项 -v, 我们可以将外部值 (并非来自stdin) 传递给awk:

```
$ VAR=10000
$ echo | awk -v VARIABLE=$VAR '{ print VARIABLE }'
1
```

还有另一种灵活的方法可以将多个外部变量传递给awk，例如：

```
$ var1="Variable1" ; var2="Variable2"
$ echo | awk '{ print v1,v2 }' v1=$var1 v2=$var2
Variable1 Variable2
```

当输入来自于文件而非标准输入时，使用：

```
$ awk '{ print v1,v2 }' v1=$var1 v2=$var2 filename
```

在上面的方法中，变量之间用空格分隔，以键-值对的形式（v1=\$var1 v2=\$var2）作为awk的命名行参数紧跟在BEGIN、{}和END语句块之后。

3. 用getline读取行

通常，grep默认读取一个文件的所有行。如果只想读取某一行，可以使用getline函数。有时候，我们需要从BEGIN语句块中读取第一行。

语法：getline var

变量var就包含了特定行的内容。

如果调用不带参数的getline，我们可以用\$0、\$1和\$2访问文本行的内容。

例如：

```
$ seq 5 | awk 'BEGIN { getline; print "Read ahead first line", $0 } {
print $0 }'
Read ahead first line 1
2
3
4
5
```

4. 用样式对awk处理的行进行过滤

我们可以为需要处理的行指定一些条件，例如：

```
$ awk 'NR < 5' # 行号小于5的行
$ awk 'NR==1,NR==4' # 行号在1到5之间的行
$ awk '/linux/' # 包含样式linux的行（可以用正则表达式来指定样式）
$ awk '!/linux/' # 不包含包含样式linux的行
```

5. 设置字段分隔符

默认的字段分隔符是空格。我们也可以使用-F "delimiter"明确指定一个分隔符：

```
$ awk -F: '{ print $NF }' /etc/passwd
```

或者

```
awk 'BEGIN { FS=":" } { print $NF }' /etc/passwd
```

在BEGIN语句块中则可以用OFS="delimiter"设置输出字段的分隔符。

6. 从awk中读取命令输出

在下面的代码中，echo会生成一个空白行。变量cmdout包含命令grep root /etc/passwd的输出，然后打印包含root的行：

将command的输出读入变量output的语法如下:

```
"command" | getline output ;
```

例如:

```
$ echo | awk '{ "grep root /etc/passwd" | getline cmdout ; print cmdout }'
root:x:0:0:root:/root:/bin/bash
```

通过使用getline, 能够将外部shell命令的输出读入变量cmdout。

awk支持以文本作为索引的关联数组。

7. 在awk中使用循环

在awk中可以使用for循环, 其格式如下:

```
for(i=0;i<10;i++) { print $i ; }
```

或者

```
for(i in array) { print array[i]; }
```

awk有很多内建的字符串控制函数, 让我们认识一下其中部分函数。

- ❑ length(string): 返回字符串的长度。
- ❑ index(string, search_string): 返回search_string在字符串中出现的位置。
- ❑ split(string, array, delimiter): 用定界符生成一个字符串列表, 并将该列表存入数组。
- ❑ substr(string, start-position, end-position): 在字符串中用字符起止偏移量生成子串, 并返回该子串。
- ❑ sub(regex, replacement_str, string): 将正则表达式匹配到的第一处内容替换成replacement_str。
- ❑ gsub(regex, replacement_str, string): 和sub()类似。不过该函数会替换正则表达式匹配到的所有内容。
- ❑ match(regex, string): 检查正则表达式是否能够匹配字符串。如果能够匹配, 返回非0值; 否则, 返回0。match()有两个相关的特殊变量, 分别是RSTART和RLENGTH。变量RSTART包含正则表达式所匹配内容的起始位置, 而变量RLENGTH包含正则表达式所匹配内容的长度。

4.8 替换文本或文件中的字符串

字符串替换是常见的文本处理任务。通过利用正则表达式匹配所需要的文本, 就可以轻松完成该任务。

4.8.1 预备知识

一听到“替换”这个词, 所有系统管理员大概都会想到sed。sed是类UNIX系统下用于文本或文件内容替换的通用工具。让我们看看如何使用sed。

4.8.2 实战演练

在4.6节中已经介绍了sed的大多数用法。你可以利用下面的方式替换一个字符串或样式：

```
$ sed 's/PATTERN/replace_text/g' filename
```

或者

```
$ stdin | sed 's/PATTERN/replace_text/g'
```

也可以用双引号 (") 代替单引号 (')。当使用双引号时，我们可以在sed样式和替换字符串中指定变量。例如：

```
$ p=pattern
$ r=replaced
$ echo "line containing a pattern" | sed "s/$p/$r/g"
line containing a replaced
```

在sed中，我们也可以不使用g。

```
$ sed 's/PATTEN/replace_text/' filename
```

这条命令会替换PATTERN匹配的第一处内容。/g意为全局(global)，这就意味着它会替换文件中所有匹配PATTERN的内容。

4.8.3 补充内容

我们已经看过了用sed进行文本替换的基本用法。接下来看一看如何将替换过的文本保存到原文件中。

将替换结果保存到文件中

当文件名传递给sed时，sed会将输出写入stdout。如果不想将输出传送至stdout，而是将更改保存到文件中，那么可以使用-i选项：

```
$ sed 's/PATTERN/replacement/' -i filename
```

例如，用另一个指定的数字替换文件中的所有3位数字：

```
$ cat sed_data.txt
11 abc 111 this 9 file contains 111 11 88 numbers 0000

$ cat sed_data.txt | sed 's/\b[0-9]\{3\}\b/NUMBER/g'
11 abc NUMBER this 9 file contains NUMBER 11 88 numbers 0000
```

上面的单行命令只替换3位数字。`\b[0-9]\{3\}\b`是一个用于匹配3位数字的正则表达式。`[0-9]`表示从0到9的数字范围。`{3}`用来匹配3次之前的数字。`\{3\}`中的`\`用于赋予(和)特殊的含义。`\b`是单词边界标记。

4.8.4 参考

4.6节中讲解了sed命令。

4.9 压缩或解压缩 JavaScript

JavaScript广泛用于网站设计。在编写JavaScript代码时，出于代码可读性与方便维护方面的考虑，我们有必要使用一些空格、注释和制表符。但是大量的空格以及制表符会增加JavaScript文件的大小。随着文件体积的增加，网页载入的时间也随之延长。因此，多数专业网站为了加快页面载入，都会对JavaScript文件进行压缩。通过压缩可以大幅度降低空白字符和换行符的数量。经过压缩的JavaScript，还可以通过加入足够的空白字符和换行符进行解压缩，这样就能够恢复代码的可读性。这则攻略就尝试在shell中发掘出类似的功能。

4.9.1 预备知识

我们准备写一个JavaScript压缩工具或代码混乱器，当然，还包括与之对应的解压缩工具。我们打算用文本与字符替换工具tr以及sed来试一试。下面来看看具体的做法。

4.9.2 工作原理

首先梳理一下对JavaScript进行压缩与解压缩的处理逻辑以及代码。

```
$ cat sample.js
function sign_out()
{

    $("#loading").show();
    $.get("log_in", {logout: "True"},

    function() {

        window.location = "";

    });

}
```

下面列出了为压缩JavaScript所需完成的工作。

- (1) 移除换行符和制表符。
- (2) 压缩空格。
- (3) 替换注释 /* 内容 */。
- (4) 替换下列内容：

- " { " 替换为 "{"
- " }" 替换为 "}"
- " (" 替换为 "("
- ")" 替换为 ")"
- " , " 替换为 ","
- " ; " 替换为 ";" (我们需要移除所有多余的空格)



要解压缩或者恢复JavaScript的可读性，我们则需要执行下面的任务：

- (1) 用 ";"\n" 替换 ";"。
- (2) 用 "{\n" 替换 "{", "\n" 替换 "}"。

4.9.3 工作原理

通过执行下面的步骤压缩JavaScript：

- (1) 移除 '\n' 和 '\t'：

```
tr -d '\n\t'
```

- (2) 移除多余的空格：

```
tr -s ' '
```

或者

```
sed 's/[ ]+//g'
```

- (3) 移除注释：

```
sed 's:/\*.*\*/::g'
```

- 因为我们需要使用 /* 和 */，所以以冒号作为sed的定界符，这样就不必对 / 进行转义了。
- * 在 sed 中被转义为 *。
- .* 用来匹配 /* 与 */ 之间所有的文本。

- (4) 移除 (、)、(、)、;、: 以及逗号前后的所有空格。

```
sed 's/ \? \([{}(){};,:]\) \? /\1/g'
```

上面的sed语句含义如下所示。

- sed代码中的 / \? \([{}(){};,:]\) \? /\1/g 用于匹配，/\1/g 用于替换。
- \([{}(){};,:]\) 用于匹配集合 [{ } () ; , :] (出于可读性方面的考虑，在这里加入了空格) 中的任意一个字符。\[(和\]) 是分组操作符，用于记忆所匹配的内容，以便在替换部分中进行向后引用。对 (和\]) 转义之后，它们便具备了另一种特殊的含义，进而可以将它们作为分组操作符。位于分组操作符前后的 \? 用来匹配可能出现在字符集合前后的空格。
- 在命令的替换部分，匹配字符串 (也就是一个可选的空格、一个来自字符集的字符再加一个可选的空格) 被匹配字符串所替换。对于匹配字符串，替换部分使用了向后引用，并通过分组操作符() 记录了匹配的字符串内容。可以用符号 \1 向后引用分组匹配的内容。

用管道将上面的步骤按照下列方式组合起来：

```
$ cat sample.js | \
tr -d '\n\t' | tr -s ' ' \
| sed 's:/\*.*\*/::g' \
| sed 's/ \? \([{}(){};,:]\) \? /\1/g'
```

得到输出：

```
function sign_out(){$("#loading").show();$.get("log_in", {logout:"True"}, function(){window.location=""}});}
```

接下来，写一个可以将这些混乱的代码恢复正常的解压缩脚本：

```
$ cat obfuscated.txt | sed 's/;/\n/g; s/{/{\n\n/g; s/}/\n\n}/g'
```

或者

```
$ cat obfuscated.txt | sed 's/;/\n/g' | sed 's/{/{\n\n/g' | sed 's/}/\n\n}/g'
```

上面的命令中：

- `s/;/\n/g` 将 `;` 替换为 `\n`;
- `s/{/{\n\n/g` 将 `{` 替换为 `{\n\n`
- `s/}/\n\n}/g` 将 `}` 替换为 `\n\n}`

4.9.4 参考

- 2.6节讲解了`tr`命令。
- 4.6节讲解了`sed`命令。

4.10 对文件中的行、单词和字符进行迭代

4

编写不同的文本处理和文件操作脚本时，通常需要对文件中的字符、单词和行进行迭代。尽管这是一个很简单的操作，但是我们会犯一些低级的错误，使所得的输出结果与期望的大相径庭。这则攻略将帮助你学习这方面的知识。

4.10.1 预备知识

将一个简单的循环用于迭代，再加上来自`stdin`或文件的重定向，这就是对文件中的行、单词和字符进行迭代的基本方法。

4.10.2 实战演练

在这里，我们要讨论对行、单词和字符进行迭代时所需执行的三个步骤。让我们来看看各个步骤是如何完成的。

(1) 迭代文件中的每一行

我们可以用一个`while`循环从标准输入中读取，因此，它在每一次迭代中都会读取一行。

下面的方法可以将`stdin`重定向到文件：

```
while read line;
do
echo $line;
done < file.txt
```

使用子`shell`的方法如下：

```
cat file.txt | { while read line; do echo $line; done }
```

命令中的cat file.txt可以用任何命令序列的输出来替换。

(2) 迭代一行中的每一个单词

我们用一个while循环迭代各行中的单词：

```
for word in $line;
do
echo $word;
done
```

(3) 迭代一个单词中的每一个字符

我们利用一个for循环对变量 i 进行迭代，迭代范围从0到字符串的长度。在每次迭代中，可以用一个特殊的记法 \${string:start_position:No_of_characters}从字符串中提取一个字符。

```
for((i=0;i<${#word};i++))
do
echo ${word:i:1} ;
done
```

4.10.3 工作原理

从文件中读取行或从行中读取单词，都是很直观的操作。不过，从单词中读取字符就需要点技巧了。为此，我们采用子串提取技术。

`${word:start_position:no_of_characters}` 返回变量word所包含的字符串中的一个子串。

`${#word}` 返回变量word的长度。

4.10.4 参考

- 1.14节讲解了Bash中不同类型的循环。
- 4.20节讲解了如何从字符串中提取字符。

4.11 按列合并文件

很多时候我们需要按列拼接文件，比如要将每一个文件的内容作为单独的一列。而cat命令所实现的拼接通常是按照行来进行的。

4.11.1 工作原理

可以用paste命令实现按列拼接，其语法如下：

```
$ paste file1 file2 file3 ...
```

让我们来尝试一下：



```
$ cat paste1.txt
1
2
3
4
5
$ cat paste2.txt
slynux
gnu
bash
hack
$ paste paste1.txt paste2.txt
1slynux
2gnu
3bash
4hack
5
```

默认的定义符是制表符，也可以用 `-d` 明确指定定义符，例如：

```
$ paste paste1.txt paste2.txt -d ","
1,slynux
2,gnu
3,bash
4,hack
5,
```

4

4.11.2 参考

4.4节讲解了如何从文本文件中提取数据。

4.12 打印文件或行中的第 n 个单词或列

我们可能有一个包含了多列数据的文件，不过只有其中的一小部分能派上用场。为了只打印相关的列或字段，我们得进行过滤。

4.12.1 预备知识

处理这种任务最为广泛的方法就是借助 `awk`。当然，用 `cut` 也可以。

4.12.2 实战演练

用下面的命令打印第5列：

```
$ awk '{ print $5 }' filename
```

也可以打印多列数据，并在列间插入指定的字符串。

比如，要打印当前目录下各文件的权限和文件名，可以使用：



```
$ ls -l | awk '{ print $1 : " $8 }'  
-rw-r--r-- : delimited_data.txt  
-rw-r--r-- : obfuscated.txt  
-rw-r--r-- : paste1.txt  
-rw-r--r-- : paste2.txt
```

4.12.3 参考

- 4.7节讲解了awk命令。
- 4.4节讲解了如何从文本文件中提取数据。

4.13 打印不同行或样式之间的文本

有时候,可能需要根据某些条件打印文本的某一部分,比如行号范围以及由起止样式所匹配的文本范围等。让我们来看看如何实现这些需求。

4.13.1 预备知识

我们可以用awk、grep和sed这类工具来根据条件打印某些文本区域。不过,我仍然觉得awk是最简单明了的方法。下面就看看如何使用awk完成这项任务。

4.13.2 实战演练

要打印出从M行到N行这个范围内的所有文本,使用下面的语法:

```
$ awk 'NR==M, NR==N' filename
```

也可以用stdin作为输入:

```
$ cat filename | awk 'NR==M, NR==N'
```

把M和N换成具体的数字:

```
$ seq 100 | awk 'NR==4,NR==6'  
4  
5  
6
```

要打印处于start_pattern与end_pattern之间的文本,使用下面的语法:

```
$ awk '/start_pattern/, /end_pattern/' filename
```

例如:

```
$ cat section.txt  
line with pattern1  
line with pattern2  
line with pattern3  
line end with pattern4  
line with pattern5
```



```
$ awk '/pa.*3/, /end/' section.txt
line with pattern3
line end with pattern4
```

用于awk中的样式为正则表达式。

4.13.3 参考

4.7节讲解了awk命令。

4.14 用脚本检验回文字符串

初次C语言编程实验课中的一个练习就是检验某个字符串是否为回文^①。这则攻略的目的就是为了告诉你如何解决类似的问题，即如何在文本中重复匹配之前的样式。

4.14.1 预备知识

sed命令能够记住之前匹配的子样式(sub pattern)。这被称为反向引用。我们可以借助这项功能解决回文问题。而在Bash中，这个问题的解决方法也不止一种。

4

4.14.2 工作原理

sed能够记得之前匹配的正则表达式，这样一来，我们就可以判断字符串中是否存在重复的字符。这种能够记忆并引用之前所匹配样式的能力就是所谓的反向引用。

让我们看看如何利用反向引用轻松地解决之前的回文问题。例如：

```
$ sed -n '/\(\.\)\1/p' filename
```

\(\.\)的作用是记录()中的子串。这里出现的. (点号)是sed用于匹配单个字符的正则表达式。

\1 对应()中匹配的第一处内容，\2 对应第二处匹配，因此我们就可以记录多个()所匹配的内容。()以\(\.\)的形式出现，表明(和)并不是普通的字符，而是有着特殊的含义。

这条sed语句会打印出所有匹配样式的连续的两个相同字符。

所有的回文字的结构特征如下。

- ❑ 如果字符数是偶数，那么它在结构上表现为：一个字符序列连着另一个字符相同但次序恰好相反的字符序列。^②
- ❑ 如果字符数是奇数，那么它在结构上表现为：一个字符序列连着另一个字符相同但次序恰好相反的字符序列，但是这两个序列中间共享一个相同的字符。^③

因此，要能够同时匹配这两种结构，我们在编写正则表达式时，需要使得其中有一个可选的

① 回文 (palindrome) 可以是单词、短语、数字等，其特点是从任何一个方向读来，无论从左至右，还是从右至左，都会得到同样的结果，例如acca、“Was it a rat I saw?”以及5885等。

② 例如ABBA。

③ 例如ABA，其中的字符B就是共享的字符。

字符。

能够匹配包含3个字符的回文字的sed正则表达式类似于下面的形式：

```
'/\(.\)\.1/p'
```

我们在字符序列及其逆序列之间加了一个额外的字符 (.)。

下面让我们来编写一个能够匹配任意长度的回文字脚本：

```
#!/bin/bash
#文件名: match_palindrome.sh
#用途: 找出给定文件中的回文字

if [ $# -ne 2 ];
then
echo "Usage: $0 filename string_length"
exit -1
fi

filename=$1 ;

basepattern='^\(.\) '

count=$(( $2 / 2 ))

for ((i=1;i<$count;i++))
do
basepattern=$basepattern'\(.\) ' ;
done

if [ $(( $2 % 2 )) -ne 0 ];
then
basepattern=$basepattern'. ' ;
fi

for ((count;count>0;count--))
do
basepattern=$basepattern'\ '$count' ;
done

basepattern=$basepattern'$/p'
sed -n "$basepattern" $filename
```

将字典文件作为输入获取一个给定长度的回文字列表，例如：

```
$ ./match_palindrome.sh /usr/share/dict/british-english 4
noon
peep
poop
sees
```

4.14.3 工作原理

这些脚本的工作过程很简单。大多数的的工作都是用来生成用于sed的正则表达式以及反向引用。

不妨通过一些例子来看看这个脚本是如何工作的。

- 如果你想匹配一个字符，并对它进行反向引用，可以用`\(.\)`匹配单个字符，再用`\1`引用它。因此，要匹配一个双字符的回文并打印，我们使用：

```
sed '/\(.\)\1/p'
```

为了指明必须从行首开始进行匹配，我们加入了行首标记`^`，于是，上面的命令就变成了`sed '/^\(.\)\1/p'`。其中，`/p`用于打印匹配内容。

- 如果我们要匹配4个字符的回文，则使用：

```
sed '/^\(.\)\(.\)\2\1/p'
```

我们用了2个`\(.\)`匹配并记录回文中前两个字符。`sed`会记住所有位于`\(`和`\)`中的匹配内容并能够反向引用它们。`\2\1`用来对所匹配的字符以相反的顺序进行反向引用。

在上面的脚本中，我们用到了一个被称为`basepattern`的变量，它包含了用于`sed`脚本的部分内容。

根据回文中字符的数量，我们利用一个`for`循环生成匹配样式。

`basepattern`最初被初始化为`basepattern="/^\(.\)'`，它对应匹配一个单字符。`for`循环用来拼接`\(.\)`和`basepattern`，循环次数为回文长度的一半。然后再用一个`for`循环以逆序的形式将反向引用拼接起来（比如`'\4\3\2\1'`），循环次数为回文长度的一半。最后，为了支持奇数长度的回文，在正则表达式与反向引用之间加上一个可选的字符`(.)`。

按照这种方法，就可以创建出用于匹配回文的`sed`匹配样式。这个匹配样式能够用来从字典文件中找出回文。

在先前的脚本中，我们用`for`循环生成`sed`命令的匹配样式。其实这个样式没必要单独生成，`sed`可以利用标签和`goto`自行实现循环结构。作为一门非常给力的语言，`sed`只需要使用一个复杂的单行脚本，就可以完成回文检查的任务。如果要从头解释就说来话长了。先动手试试下面的脚本：

```
$ word="malayalam"

$ echo $word | sed ':loop ; s/^\(.\)(.*)\1\2/; t loop; /^.\?$/ { s/.*/PALINDROME/; q; }; s/.*/NOT PALINDROME/ '
```

PALINDROME

如果你对`sed`脚本非常感兴趣，不妨看看一本专门讲解`sed`和`awk`的参考书，即由Dale Dougherty和Arnold Robbins所著的《`sed`与`awk`（第2版）》。

你可以利用这本书试着分析一下上面那个测试回文的单行`sed`脚本。

4.14.4 补充内容

下面来看看与该任务相关的其他选项或是内容。

最简单直接的方法

检查一个字符串是否是回文的最简单的方法就是使用`rev`命令。

rev命令接受一个文件或stdin作为输入，并逆序打印出每一行内容。

试试下面的代码：

```
string="malayalam"
if [[ "$string" == "$(echo $string | rev )" ]];
then
echo "Palindrome"
else
echo "Not palindrome"
fi
```

rev命令可以结合其他命令来解决不同的问题。来看一个将句中所有单词顺序反转的例子：

```
sentence='this is line from sentence'
echo $sentence | rev | tr ' ' '\n' | tac | tr '\n' ' ' | rev
```

输出如下：

```
sentence from line is this
```

在上面的单行命令中，首先用rev命令反转所有的字符，然后用tr命令将空格替换成 \n，这样就将所有的单词分隔成每行一个，接着通过tac命令将所有的行进行反转，再用tr将反转后的行合并成单行。最后，对合并后的单行使用rev，实现单词反转。

4.14.5 参考

- 4.6节讲解了sed命令。
- 1.15节讲解了字符串比较操作符。

4.15 以逆序形式打印行

这一节的内容很简单。这里讲的东西可能看起来没什么用，不过它可以用来在Bash中模拟栈结构。所以这节还是有些意思的。下面我们要将一个文件中的文本行以逆序形式打印出来。

4.15.1 预备知识

只需要用点awk的小技巧就能完成这项任务。不过，命令tac可以直接做到同样的事。这个命令的名称其实就是反过来书写的cat。

4.15.2 实战演练

先来试试tac。该命令的语法如下：

```
tac file1 file2 ...
```

它也可以从stdin中读取：

```
$ seq 5 | tac
5
```

```
4
3
2
1
```

在tac中, \n是默认的行分隔符。但我们也可以用 -s "分隔符"选项指定自己的分隔符。使用awk的实现方式如下:

```
$ seq 9 | \
awk '{ lifo[NR]=$0; lno=NR }
END{ for(;lno>-1;lno--){ print lifo[lno]; }
}'
```

在shell脚本中, \可以很方便地将单行命令分解成多行。

4.15.3 工作原理

这个awk脚本非常简单。我们将每一行都存入一个关联数组中, 用行号作为数组索引(行号由NR给出), 最后由awk执行END语句块。为了得到最后一行的行号, 在{}语句块中使用lno=NR。因此, 这个脚本从最后一行一直迭代到第0行, 将存储在数组中的各行以逆序方式打印出来。

4.15.4 参考

4.19节讲解了如何使用awk编写tac。

4.16 解析文本中的电子邮件地址和 URL

从给定的文件中解析出所需要的文本是我们从事文本处理时常见的一项任务。诸如电子邮件地址、URL等都能够借助适合的正则表达式找出来。在大多数情况下, 我们需要从一个包含大量无关字符及单词的电子邮件客户列表或HTML网页中将电子邮件地址解析提取出来。

4.16.1 预备知识

这个问题可以用egrep解决。

4.16.2 实战演练

能够匹配一个电子邮件地址的egrep正则表达式如下:

```
[A-Za-z0-9.]+@[A-Za-z0-9.]+\.[a-zA-Z]{2,4}
```

例如:

```
$ cat url_email.txt
```

```
this is a line of text contains, <email> #slynux@slynux.com. </email>
and email address, blog "http://www.google.com", test@yahoo.com
```

```

dfdfdfdfdfdf;cool.hacks@gmail.com<br />

<ahref="http://code.google.com"><h1>Heading</h1>

$ egrep -o '[A-Za-z0-9.]+@[A-Za-z0-9.]+\.[a-zA-Z]{2,4}' url_email.txt
slynux@slynux.com
test@yahoo.com
cool.hacks@gmail.com

匹配一个HTTP URL的egrep正则表达式如下:

http://[a-zA-Z0-9\-.]+\.[a-zA-Z]{2,4}

例如:

$ egrep -o "http://[a-zA-Z0-9.]+\.[a-zA-Z]{2,3}" url_email.txt
http://www.google.com
http://code.google.com

```

4.16.3 工作原理

如果逐个部分进行设计的话,这些正则表达式其实很简单。在匹配电子邮件地址的正则表达式中,我们都知道电子邮件地址可以采用name@domain.some_2-4_letter这样的形式。那么,在编写正则表达式时,也要遵循同样的规则:

```
[A-Za-z0-9.]+@[A-Za-z0-9.]+\.[a-zA-Z]{2,4}
```

[A-Za-z0-9.]+ 表示在表示字面意义的字符@出现之前,[]中的字符组合应该出现一次或多次(这也正是+的含义)。然后[A-Za-z0-9.]同样应该出现一次或多次(+).样式\表示应该呈现一个表示字面意义的“.”(点号),而[a-zA-Z]{2,4}表示字母的长度应该在2到4之间(包括2和4)。

匹配HTTP URL与匹配电子邮件地址类似,只是不需要匹配name@部分。

```
http://[a-zA-Z0-9.]+\.[a-zA-Z]{2,3}
```

4.16.4 参考

- 4.6节讲解了sed命令。
- 4.2节讲解了如何使用正则表达式。

4.17 打印文件中某个样式之前或之后的n行

在文本处理中经常需要打印由样式匹配的某个文本区域。有时候,我们也许需要位于样式之前或之后的若干行文本。例如有一个文件包含了电影演员的评级情况,文件中的每一行都对一位电影演员的详细信息,而我们需要从中找出某位演员的评级以及最接近这位演员的其他演员的详细信息。来看看如何实现这个需求。

4.17.1 预备知识

grep是在文件中搜索文本的最佳工具。通常，grep会打印出匹配给定样式的文本行或文本。不过grep的前后行（context line）控制选项使它可以打印出位于样式匹配行之前、之后的或者同时包含前后的文本行。

4.17.2 实战演练

这个技巧最好通过一个电影演员名单来讲解。例如：

```
$ cat actress_rankings.txt | head -n 20
1 Keira Knightley
2 Natalie Portman
3 Monica Bellucci
4 Bonnie Hunt
5 Cameron Diaz
6 Annie Potts
7 Liv Tyler
8 Julie Andrews
9 Lindsay Lohan
10 Catherine Zeta-Jones
11 CateBlanchett
12 Sarah Michelle Gellar
13 Carrie Fisher
14 Shannon Elizabeth
15 Julia Roberts
16 Sally Field
17 TéaLeoni
18 Kirsten Dunst
19 Rene Russo
20 JadaPinkett
```

要打印出匹配“Cameron Diaz”的文本行及其之后的3行，使用下面的命令：

```
$ grep -A 3 "Cameron Diaz" actress_rankings.txt
5 Cameron Diaz
6 Annie Potts
7 Liv Tyler
8 Julie Andrews
```

要打印出匹配行及其之前的3行，使用下面的命令：

```
$ grep -B 3 "Cameron Diaz" actress_rankings.txt
2 Natalie Portman
3 Monica Bellucci
4 Bonnie Hunt
5 Cameron Diaz
```

打印出匹配行及其之前的2行和之后的2行，使用下面的命令：

```
$ grep -C 2 "Cameron Diaz" actress_rankings.txt
3 Monica Bellucci
4 Bonnie Hunt
```



5 Cameron Diaz
6 Annie Potts
7 Liv Tyler

你是否还在想我是从哪里搞到这些评级信息的？

我可是曾经只用基本的sed、awk和grep命令就解析了一个包含大量图片和HTML内容的网站。详情请参阅第5章。

4.17.3 参考

4.3节讲解了grep命令。

4.18 在文件中移除包含某个单词的句子

只要能写出正确的正则表达式，移除包含某个单词的句子简直就是手到擒来。这里给出了一个解决类似问题的练习。

4.18.1 预备知识

sed是进行文本替换的不二之选。这样，我们就可以通过sed用空白替代匹配的句子。

4.18.2 实战演练

先创建一个包含替换文本的文件。例如：

```
$ cat sentence.txt
```

```
Linux refers to the family of Unix-like computer operating systems  
that use the Linux kernel. Linux can be installed on a wide variety of  
computer hardware, ranging from mobile phones, tablet computers and video  
game consoles, to mainframes and supercomputers. Linux is predominantly  
known for its use in servers. It has a server market share ranging  
between 20-40%. Most desktop computers run either Microsoft Windows or  
Mac OS X, with Linux having anywhere from a low of an estimated 1-2% of  
the desktop market to a high of an estimated 4.8%. However, desktop use  
of Linux has become increasingly popular in recent years, partly owing  
to the popular Ubuntu, Fedora, Mint, and openSUSE distributions and the  
emergence of netbooks and smart phones running an embedded Linux.
```

我们的目标是移除包含单词“mobile phones”的句子。用下面的sed语句来完成这项任务：

```
$ sed 's/ [^]*mobile phones[^]*\./g' sentence.txt
```

```
Linux refers to the family of Unix-like computer operating systems  
that use the Linux kernel. Linux is predominantly known for its use  
in servers. It has a server market share ranging between 20-40%. Most  
desktop computers run either Microsoft Windows or Mac OS X, with Linux  
having anywhere from a low of an estimated 1-2% of the desktop market to  
a high of an estimated 4.8%. However, desktop use of Linux has become
```

increasingly popular in recent years, partly owing to the popular Ubuntu, Fedora, Mint, and openSUSE distributions and the emergence of netbooks and smart phones running an embedded Linux.

4.18.3 工作原理

让我们分析一下sed的正则表达式's/ [^.]*mobile phones[^.]*\./g'。

该正则表达式的格式为：'s/匹配样本/替代字符串/g'。

它将与匹配样本相匹配的每一处内容都用替代字符串进行替换。

这里的匹配样本是用来匹配一句文本的正则表达式。文件中的每一句话第一个字符都是空格，句与句之间都以"."来分隔。因此我们需要匹配内容的格式就是：空格+若干文本+需要匹配的字符串+若干文本+句点。一个句子除了作为定界符的句点之外，可以包含任意字符。因此我们要使用[^.]. [^.]* 可以匹配除句点之外的任何字符的组合。用来匹配文本的“mobile”被放置在两个[^.]*之间。每一个匹配的句子均被//替换（注意，/与/之间没有任何内容）。

4.18.4 参考

□ 4.6节讲解了sed命令。

□ 4.2节讲解了如何使用正则表达式。

4

4.19 用 awk 实现 head、tail 和 tac

只有通过实践才能够掌握文本处理的各种操作。这则攻略将帮助你综合运用之前已学过的一些命令。

4.19.1 预备知识

命令head、tail、uniq和tac都是逐行操作的。不管什么时候，只要我们需要进行逐行处理，我们都可以用awk来解决。接下来让我们利用awk模拟这些命令。

4.19.2 实战演练

看一看如何用不同的awk命令来模拟诸如head、tail、tac等命令。

模拟head命令读取文件前10行并打印出来：

```
$ awk 'NR <=10' filename
```

模拟tail命令打印文件的后10行：

```
$ awk '{ buffer[NR % 10] = $0; } END { for(i=1;i<11;i++) { print buffer[i%10] } }' filename
```

模拟tac命令逆序打印输入文件的所有行：

```
$ awk '{ buffer[NR] = $0; } END { for(i=NR; i>0; i--) { print buffer[i] } }' filename
```

4.19.3 工作原理

在head的awk实现中，我们打印输入流中行号小于或等于0的行。行号可以通过特殊变量NR获得。

在tail命令的实现中，我们使用了散列技术。数组buffer的索引是由散列函数NR%10决定的，其中变量NR包含了当前行的行号，\$0包含了当前的文本行。因此，% 将所有余数相同的行都映射到一个特定的数组索引中。在END()语句块中，对数组的10个索引进行迭代，并打印出对应的行。

在tac命令的实现中，只是将所有的行存入一个数组中。当程序流程进入END()语句块时，NR存储着最后一行的行号。然后在for循环中对NR递减到1，并在每一次循环中打印存储的行。

4.19.4 参考

- 4.7节讲解了awk命令。
- 3.14节讲解了uniq命令。
- 2.8节讲解了uniq命令。
- 4.15节讲解了tac命令。

4.20 文本切片与参数操作

这则攻略考查了一些简单的文本替换技术以及Bash中可用的参数扩展简写法。这些简单的技巧通常能够让我们免于敲入多行代码之苦。

4.20.1 实战演练

就用下面的任务练练手吧。

替换变量内容中的部分文本：

```
$ var="This is a line of text"
$ echo ${var/line/REPLACED}
This is a REPLACED of text"
```

line被REPLACED替换。

我们可以通过指定字符串的起始位置和长度来生成子串，语法如下：

```
${variable_name:start_position:length}
```

用下面的命令可以打印第5个字符之后的内容：

```
$ string=abcdefghijklmnpqrstuvwxyx
$ echo ${string:4}
efghijklmnpqrstuvwxyx
```



从第5个字符开始，打印8个字符：

```
$ echo ${string:4:8}  
efghijkl
```

起始字符的索引从0开始计数。我们也可以从后向前计数，将最后一个字符索引记为-1。但如果使用负数作为索引值的话，必须将负数放入括号内，例如(-1)就是最后一个字符的索引。

```
echo ${string: -1}  
z  
$ echo ${string: -2: 2}  
yz
```

4.20.2 参考

4.10节讲解了从单词中切分字符。



本章内容

- 网站下载
- 以格式化纯文本形式下载网页
- cURL入门
- 从命令行访问Gmail
- 解析网站数据
- 制作图片抓取器及下载工具
- 创建网页相册生成器
- 建立一个Twitter命令行客户端
- 基于Web后端的定义查询工具
- 查找网站中的无效链接
- 跟踪网站变更
- 以POST方式发送网页并读取响应

5.1 入门

Web正在成为反映技术发展的晴雨表，数据处理都要经由它来完成。尽管shell脚本没法像PHP一样在Web上大包大揽，但总还是有一些活儿适合它。本章会研究一些用于解析网站内容、下载数据、发送数据表单以及网站任务自动化之类的实例。我们可以仅用几行脚本就将很多原本需要通过浏览器交互进行的活动管理自动化。通过命令行工具利用HTTP协议所提供的功能，我们可以用脚本解决大部分Web自动化的问题。来尽情享受接下来的内容吧。

5.2 网站下载

从给定的URL中下载文件或网页很简单，一些命令行下载工具就可以完成这项任务。

5.2.1 预备知识

wget是一个用于文件下载的命令行工具，选项多且用法活。

5.2.2 实战演练

用wget可以下载网页或远程文件：

```
$ wget URL
```

例如:

```
$ wget http://slynuux.org
--2010-08-01 07:51:20-- http://slynuux.org/
Resolving slynuux.org... 174.37.207.60
Connecting to slynuux.org[174.37.207.60]:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 15280 (15K) [text/html]
Saving to: "index.html"

100%[=====>] 15,280 75.3K/s in
0.2s

2010-08-01 07:51:21 (75.3 KB/s) - "index.html" saved [15280/15280]
```

可以指定从多个URL处进行下载:

```
$ wget URL1 URL2 URL3..
```

可以依据URL用wget下载文件:

```
$ wget ftp://example_domain.com/somefile.img
```

通常下载的文件名和URL中的文件名保持一致, 下载信息或进度被写入stdout。

你可以通过选项-o指定输出文件名^①。如果存在同名文件, 那么会先将该同名文件清空(truncate)再将下载文件写入。

你也可以用选项-o指定一个日志文件, 从而不必将日志信息打印到stdout。

```
$ wget ftp://example_domain.com/somefile.img -O dloaded_file.img -o log
```

运行该命令, 屏幕上不会出现任何内容。日志或进度信息被写入文件log, 输出文件为dloaded_file.img。

由于不稳定的Internet连接, 下载有可能被迫中断。我们可以将重试次数作为命令参数, 这样一旦下载中断, wget在放弃下载之前还会继续进行多次尝试。

用-t指定重试次数:

```
$ wget -t 5 URL
```

5.2.3 补充内容

wget还有一些其他选项, 用以应对不同的情况。来看看其中的一部分吧。

1. 下载限速

当我们的下载带宽有限, 却又有多个应用程序共享该Internet连接时, 进行大文件下载往往将榨干所有的带宽, 严重阻滞其他进程。wget命令有一个内建的选项可以限定下载任务能够占有的最大带宽, 从而使得其他应用程序流畅运行。

^① 选项'-O file'中的-o并不是表示不采用URL中的文件名, 而是由file指定。它的作用类似于shell重定向, 也就是说, 命令wget -O file http://foo实际上相当于wget -O - http://foo > file。——译者注

我们可以用`--limit-rate`按照下面的方式对`wget`限速：

```
$ wget --limit-rate 20k http://example.com/file.iso
```

在命令中用`k`（千字节）和`m`（兆字节）指定速度限制。

还可以指定最大下载配额（quota）。配额一旦用尽，下载随之停止。在下载多个文件的时候，对总下载量进行限制是很必要的，这能够避免在无意中占用过多磁盘空间。

使用`--quota`或`-Q`：

```
$ wget -Q 100m http://example.com/file1 http://example.com/file2
```

2. 断点续传

如果使用`wget`进行的下载在完成之前被中断，可以利用选项`-c`从断点开始继续下载：

```
$ wget -c URL
```

3. 用cURL下载

`cURL`是另一个高级命令行工具，功能比`wget`更胜一筹。

可以按照下面的方式用`cURL`下载：

```
$ curl http://slynux.org > index.html
```

和`wget`不同，`curl`并不将下载数据写入文件，而是写入标准输出（`stdout`）。因此我们必须用重定向操作符将数据从`stdout`重定向到文件。

4. 复制或镜像整个网站

`wget`有一个选项可以使其像爬虫一样以递归的方式收集网页上所有的URL链接，并逐个下载，这样一来，我们就能够下载一个网站的所有页面。

要实现这个任务，可以按照下面的方式使用选项`--mirror`：

```
$ wget --mirror exampledomain.com
```

或者

```
$ wget -r -N -l DEPTH URL
```

`-l`指定页面层级`DEPTH`，这意味着`wget`只会向下遍历指定的页面级数。这个选项要与`-r`（recursive，递归选项）一同使用。另外，`-N`允许对文件使用时间戳，URL表示下载的网站起始地址。

5. 访问需要认证的HTTP或FTP页面

一些网页需要HTTP或FTP认证，可以用`--user`和`--password`提供认证信息：

```
$ wget --user username --password pass URL
```

也可以不在命令行中指定密码，而由网页提示并手动输入密码，这就需要将`--password`改为`--ask-password`。

5.3 以格式化纯文本形式下载网页

网页其实就是包含HTML标记和其他诸如JavaScript、CSS等元素的HTML页面。HTML标记是网页的基础。我们也许需要解析网页来查找特定的内容，这时Bash就能派上用场了。当我们下

载网页的时候，实际接收到的就是一个HTML文件。这种格式化的数据需要用浏览器查看。但在大多数情况下，解析一个格式化文本文件要比解析HTML数据来得容易。因此如果我们能够得到一个格式类似于浏览器页面那样的文本文件，就能省下不少为剥离HTML标签而花费的功夫。Lynx是一个有意思的基于命令行的Web浏览器。我们实际上可以将Lynx的纯文本格式化输出作为网页来获取。来看看这一切该如何实现。

实战演练

用lynx命令的-dump选项将网页以ASCII字符的形式下载到文本文件中：

```
$ lynx -dump URL > webpage_as_text.txt
```

这个命令会将所有的超链接()作为文本输出的页脚列在References标题之下。这就省得我们再用正则表达式单独解析链接。

例如：

```
$ lynx -dump http://google.com > plain_text_page.txt
```

你可以用cat命令查看纯文本版的网页：

```
$ cat plain_text_page.txt
```

5.4 cURL 入门

5

作为一款强力工具，cURL支持包括HTTP、HTTPS、FTP在内的众多协议。它还支持POST、cookie、认证、从指定偏移处下载部分文件、参照页(referer)、用户代理字符串、扩展头部(extra header)、限速、文件大小限制、进度条等特性。如果你想将网页处理流程及数据检索自动化，那么cURL会助你一臂之力。这则攻略将为你展示cURL一系列最为重要的特性。

5.4.1 预备知识

在默认情况下，主流Linux发行版中并没有包含cURL，你得使用包管理器进行安装。不过多数发行版都默认附带了wget。

cURL通常将下载文件输出到stdout，将进度信息输出到stderr。要想避免显示进度信息，请使用--silent选项。

5.4.2 实战演练

curl命令可以用来执行下载、发送各种HTTP请求、指定HTTP头部等操作。让我们看看用cURL如何实现这些任务。

```
$ curl URL --silent
```

上面的命令将下载文件输出到终端（所下载的数据都被写入stdout）。

--silent选项使得curl命令不显示进度信息。如果需要这些信息，将--silent移除即可。

```
$ curl URL --silent -O
```

选项 -O 用来将下载数据写入文件，而非写入标准输出。该文件采用的是从URL中解析出的文件名。

例如：

```
$ curl http://slynx.org/index.html --silent -O
```

这将创建文件index.html。

这条命令不再将网页写入stdout，而是写入和URL中相同文件名的文件中。因此要确保这是一个指向远程文件的URL。curl http://slynx.org -O --silent将会显示错误信息，因为无法从URL中解析出文件名。

```
$ curl URL --silent -o new_filename
```

选项 -o 用来将下载的数据写入指定名称的文件中。

如果需要在下载过程中显示形如 # 的进度条，用 --progress 代替 --silent。

```
$ curl http://slynx.org -o index.html --progress
##### 100.0%
```

5.4.3 补充内容

在5.3节中，我们学习了如何下载文件以及如何将HTML页面打印到终端。cURL还包括一些高级选项，让我们进一步研究一下。

1. 断点续传

和wget不同，cURL包含更高级的下载恢复特性，能够从特定的文件偏移处继续下载。它可以通过指定一个偏移量来下载部分文件。

```
$ curl URL/file -C offset
```

偏移量是以字节为单位的整数。

如果我们只是想断点续传，那么cURL不需要我们知道准确的字节偏移。要是你希望cURL推断出正确的续传位置，请使用选项 -C -，就像这样：

```
$ curl -C - URL
```

cURL会自动计算出应该从哪里开始续传。

2. 用cURL设置参照页字符串

参照页（referer^①）是位于HTTP头部中的一个字符串，用来标识用户是从哪个页面到达当前页面的。如果用户点击了网页A中的某个链接，那么用户就会转到网页B，网页B头部的参照页字符串会包含网页A的URL。

^① referer其实应该是英文单词referrer，不过拼错的人太多了，所以编写标准的人也就将错就错了。

网页可以根据条件进行判断, 如果参考页是www.google.com, 那么就返回一个Google页面, 否则返回其他页面。

你可以用curl命令的--referer选项指定参考页字符串:

```
$ curl --referer Referer_URL target_URL
```

例如:

```
$ curl --referer http://google.com http://slynux.org
```

3. 用cURL设置cookie

我们可以用curl来存储HTTP操作过程中使用到的cookie。

要指定cookie, 使用--cookie "COOKIES"选项。

COOKIES需要以name=value的形式来给出。如果要指定多个cookie, 用分号分隔, 例如:

```
$ curl http://example.com --cookie "user=slynux;pass=hack"
```

如果要将cookie另存为一个文件, 使用--cookie-jar选项。例如:

```
$ curl URL --cookie-jar cookie_file
```

4. 用cURL设置用户代理字符串

如果不指定用户代理, 一些需要检验用户代理的网页就无法显示。你可能已经注意到有些网站只能在Internet Explorer (IE) 下正常工作, 如果使用其他浏览器, 这些网站就会提示说它只能用IE访问。这是因为这些网站检查了用户代理。你可以用curl把用户代理设置成IE, 而后你就会发现网页又能正常显示了。

用cURL的--user-agent或-A选项可以设置用户代理:

```
$ curl URL --user-agent "Mozilla/5.0"
```

其他HTTP头部信息也可以通过cURL来发送。用-H"头部信息"传递多个头部信息。例如:

```
$ curl -H "Host: www.slynux.org" -H "Accept-language: en" URL
```

5. 限定cURL可占用的带宽

如果带宽有限, 又有多个用户共享, 为了平稳流畅地分享带宽, 我们可以用--limit-rate限制curl的下载速度:

```
$ curl URL --limit-rate 20k
```

在命令中用k (千字节) 和m (兆字节) 指定下载速度限制。

6. 指定最大下载量

可以用--max-filesize选项指定可下载的最大文件大小:

```
$ curl URL --max-filesize bytes
```

如果文件大小超出限制, 命令返回一个非0的退出码。如果命令正常运行, 返回0。

7. 用cURL进行认证

可以用cURL的选项-u完成HTTP或FTP认证。

-u username:password可用来指定用户名和密码。它也可以不指定密码, 而在后续的运

行过程中经提示输入密码。

如果你喜欢经提示后输入密码，只需要使用 `-u username` 即可。例如：

```
$ curl -u user:pass http://test_auth.com
```

要使用密码提示的话，则用：

```
$ curl -u user http://test_auth.com
```

8. 只打印响应头部信息（不包括数据部分）

只打印响应头部（response header）有助于进行各种检查或统计。例如，要求无须下载整个页面内容就能够检验某个页面是否能够打开，那么我们只用读取HTTP响应头部就能够知道这个页面是否可用。

检查HTTP头部的一个用例就是在下载之前先查看文件大小。我们可以在下载之前，通过检查HTTP头部中的 `Content-Length` 参数来得知文件的长度。同样还可以从头部检索出其他一些有用的参数。`Last-Modified` 参数能告诉我们远程文件最后的改动时间。

通过 `-I` 或 `-head` 就可以只打印HTTP头部信息，而无须下载远程文件。例如：

```
$ curl -I http://slynx.org
HTTP/1.1 200 OK
Date: Sun, 01 Aug 2010 05:08:09 GMT
Server: Apache/1.3.42 (Unix) mod_gzip/1.3.26.1a mod_log_bytes/1.2 mod_bwlimited/1.4
mod_auth_passthrough/1.8 FrontPage/5.0.2.2635 mod_ssl/2.8.31 OpenSSL/0.9.7a
Last-Modified: Thu, 19 Jul 2007 09:00:58 GMT
ETag: "17787f3-3bb0-469f284a"
Accept-Ranges: bytes
Content-Length: 15280
Connection: close
Content-Type: text/html
```

5.4.4 参考

5.13节

5.5 从命令行访问 Gmail

Gmail (<http://mail.google.com>) 是 Google 被广泛使用的一项免费电子邮件服务。你可以使用经过认证的 RSS feed 来读取个人邮件。我们也能够通过发送人姓名以及邮件主题来解析 RSS feed。这使得我们无须打开网页浏览器就能够查看收件箱中的未读邮件。

5.5.1 实战演练

来看下面这个脚本文件，它的作用是通过解析 Gmail 的 RSS feed 来显示未读的邮件：

```
#!/bin/bash
文件名: fetch_gmail.sh
#用途: 获取 Gmail 工具
```

```

username="PUT_USERNAME_HERE"
password="PUT_PASSWORD_HERE"

SHOW_COUNT=5 # 最近未读的邮件数目

echo

curl -u $username:$password --silent "https://mail.google.com/mail/
feed/atom" | \
tr -d '\n' | sed 's:</entry>:\n:g' | \
sed 's/.*<title>\(.*\)</title>.*<author><name>\([^<]*\)</name><email>
\([^<]*\).*/Author: \2 [\3] \nSubject: \1\n/' | \
head -n ${SHOW_COUNT * 3}

```

输出如下:

```

$ ./fetch_gmail.sh
Author: SLYNUX [ slynux@slynux.com ]
Subject: Book release - 2

Author: SLYNUX [ slynux@slynux.com ]
Subject: Book release - 1
.
.....共5封未读邮件

```

5.5.2 工作原理

这个脚本通过用户认证使用cURL来下载RSS feed。用户认证信息由-u username: password 参数提供。也可以只使用-u user, 这样就不用再提供密码。然后由cURL在运行时以交互的方式索取密码。

我们将管道命令拆分成不同的部分来解释这个脚本的工作原理。

tr -d '\n' 移除了所有的换行符, 这样我们就可以将\n作为定界符来重建邮件项。sed 's:</entry>:\n:g' 将每一处</entry> 替换成换行符, 以保证每一条邮件项各自单独为行, 从而就能够逐条解析邮件。通过查看<https://mail.google.com/mail/feed/atom>页面源代码中用于RSS feed的XML标记, 我们可以看出每一条邮件项都对应对应着<entry> TAGS </entry>。

该脚本接下来的部分如下:

```

sed 's/.*<title>\(.*\)</title>.*<author><name>\([^<]*\)</
name><email>
\([^<]*\).*/Author: \2 [\3] \nSubject: \1\n/'

```

脚本用<title>\(.*\)</title>匹配邮件标题, 用<author><name>\([^<]*\)</ name>匹配发件人姓名, <email>\([^<]*\)</email>匹配发件人电子邮件地址。接下来, 按照下面的方法使用反向引用。

- Author: \2 [\3] \nSubject: \1\n用来将匹配内容以易于阅读的格式来替换原先的邮件项。\\1对应于第一处匹配, \\2对应于第二处匹配, 依次类推。
- SHOW_COUNT=5用来设置需要在终端中显示的未读邮件数量。

- head用来显示SHOW_COUNT*3行文本^①。SHOW_COUNT乘以3是因为每一封未读邮件的相关信息需要占用3行。

5.5.3 参考

- 5.4节讲解了curl命令。
- 4.6节讲解了sed命令。

5.6 解析网站数据

消除不必要的细节有助于解析网页数据。sed和awk是完成这项工作的主要工具。在4.3节中，我们见到过一份演员评级列表。这个列表就是通过解析<http://www.johntorres.net/BoxOfficefemaleList.html>得到的。

让我们看看如何使用文本处理工具来解析同样的数据。

5.6.1 实战演练

请看下面用于从网站解析演员详细信息的命令序列：

```
$ lynx -dump http://www.johntorres.net/BoxOfficefemaleList.html | \
grep -o "Rank-.*" | \
sed 's/Rank-//; s/\[[0-9]\+\]//; ' | \
sort -nk 1 | \
awk '
{
  for(i=3;i<=NF;i++){ $2=$2" "$i }
  printf "%-4s %s\n", $1,$2 ,
}' > actresslist.txt
```

输入如下：

```
# 由于版面空间限制，只显示前3位演员的信息
1 Keira Knightley
2 Natalie Portman
3 Monica Bellucci
```

5.6.2 工作原理

Lynx是一个基于命令行的网页浏览器。它并不会为我们显示出一堆原始的HTML代码，而是能够打印出网站的文本版本，这个文本版和我们在浏览器中看到的页面一模一样。这样一来，就免去了移除HTML标记的工作。我们按照下面的方法用sed解析以Rank开头的行：

```
sed 's/Rank-//; s/\[[0-9]\+\]//; '
```

① 这里的SHOW_COUNT*3行文本并不包括脚本开始部分由echo生成的那一行（空行）。

然后再根据评级情况对这些行排序。这里我们要使用`awk`来指定宽度，以保持演员的评级信息和姓名之间的空隙。`%-4s`指定了4个字符的宽度。除了第一个字段外，所有的字段都用`$2`拼接成了一个字符串。

5.6.3 参考

- 4.6节讲解了`sed`命令。
- 4.7节讲解了`awk`命令。
- 5.3节讲解了`lynx`命令。

5.7 制作图片抓取器及下载工具

当需要下载某个网页上所有的图片时，图片抓取器（crawler）能够帮上我们的大忙。用不着把HTML源码翻个底朝天来摘取出所有的图片，我们可以用脚本解析图像文件并将它们自动下载下来。来看看这是如何实现的。

5.7.1 实战演练

用于从网页上抓取并下载图片的Bash脚本如下：

```
#!/bin/bash
#用途:图片下载工具
#文件名:img_downloader.sh
if [ $# -ne 3 ];
then
    echo "Usage: $0 URL -d DIRECTORY"
    exit -1
fi

for i in {1..4}
do
    case $1 in
        -d) shift; directory=$1; shift ;;
        *) url=${url:-$1}; shift;;
    esac
done

mkdir -p $directory;
baseurl=$(echo $url | egrep -o "https?://[a-z.]+")

curl -s $url | egrep -o "<img src=[^>]*>" |
sed 's/ /tmp/$$.list

sed -i "s/^\/!$baseurl\/!/" /tmp/$$.list

cd $directory;
```



```
while read filename;
do
    curl -s -O "$filename" --silent

done < /tmp/$$.list
```

使用方法:

```
$ ./img_downloader.sh http://www.flickr.com/search/?q=linux -d images
```

5.7.2 工作原理

上述图片下载器脚本首先解析HTML页面,除去之外的所有标记,然后从标记中解析出src="URL"并将图片下载到指定的目录中。这个脚本接受一个网页URL和用于存放图片的目录路径作为命令行参数。脚本开始部分用了一种巧妙的方式来解析命令行参数。[\$# -ne 3] 检查脚本参数数量是否为3个,如果不是,它就会退出运行,同时返回脚本用法说明。

如果是3个参数,那么就解析URL和目标目录。解析过程需要用到一些技巧:

```
for i in {1..4}
do
    case $1 in
        -d) shift; directory=$1; shift ;;
        *) url=${url:-$1}; shift;;
    esac
done
```

用for循环执行4次迭代(第4次迭代其实也没有什么影响,无非就是执行了几次case语句)。

case语句会对第一个参数(\$1)求值,以匹配-d等被检查到的字符串参数。我们可以像下面那样将-d置于命令行的任意位置:

```
$ ./img_downloader.sh -d DIR URL
```

或者

```
$ ./img_downloader.sh URL -d DIR
```

shift用来移动参数。当使用shift后,\$2的值就被赋给\$1,如果再次使用shift,则\$3的值被赋给\$1,往后依次类推,因此我们通过\$1就可以对所有的参数进行求值。

如果匹配-d(-d),显然下一个参数就是目标目录的值。*)对应默认匹配(default match)。它能够匹配除了-d之外的任何内容。因此在默认匹配中,当\$1=""或\$1=URL时,我们需要采用\$1=URL,避免"*"将变量url覆盖掉。所以我们使用了url=\${url:-\$1}。如果url不为空,它会返回URL值,否则返回\$1的值。

egrep -o "]*>"只打印包括属性值在内的标记。[^>]*用来匹配除>之外的所有字符,结果就是。

sed 's/标记中得到所有的图像文件URL。

图像文件源路径有两种类型：相对路径和绝对路径。绝对路径包含以`http://`或`https://`起始的完整URL，相对路径则以`/`或图像文件名起始。

绝对路径：`http://example.com/image.jpg`

相对路径：`/image.jpg`

对于以起始的相对路径，应该用基址URL (base URL) 把它转换为`http://example.com/image.jpg`。

我们一开始通过解析，将基址URL存入变量`baseurl`中，以便进行路径转换。

然后，用`sed -i "s!~/!$baseurl/!" /tmp/$$.list`将所有的`/`替换成`baseurl`。

`while`循环用来对图片的URL列表进行逐行迭代，并用`curl`下载图像文件。与`curl`一同使用的`-silent`选项可避免下载进度信息出现在屏幕上。

5.7.3 参考

- 5.4节讲解了`curl`命令。
- 4.6节讲解了`sed`命令。
- 4.3节讲解了`grep`命令。

5.8 网页相册生成器

Web开发人员通常都会为网站设计相册页面，这些页面上包含着多个图像缩略图。点击缩略图，就会出现一幅放大的图片。但如果需要很多图片的话，每一次都得复制 `` 标记、调整图片大小来创建缩略图、把调整好的图片放进缩略图目录、测试链接，诸如此类的活儿实在烦琐。这些工作会花费大量时间，但都是一些简单无谓的重复。我们可以写一个脚本将这些工作步骤实现自动化处理。这样一来，创建缩略图、将缩略图放入对应的目录、生成用于 `` 标记的代码片段都可以在短短几秒钟内自动搞定。这则攻略正是要告诉你如何实现刚才所说的一切。

5

5.8.1 预备知识

要完成这项任务，可以通过`for`循环对当前目录下的所有图片进行迭代。这需要借助一些常见的Bash工具，如`cat`和`convert`。我们将用所有的图片生成一个HTML相册`index.html`。要使用`convert`，先确保你已经安装了`Imagemagick`。

5.8.2 实战演练

生成HTML相册页面的Bash脚本如下：

```
#!/bin/bash
# 文件名: generate_album.sh
# 用途: 用当前目录下的图片创建相册

echo "Creating album..."
```

```
mkdir -p thumbs
cat <<EOF > index.html
<html>
<head>
<style>
body
{
    width:470px;
    margin:auto;
    border: 1px dashed grey;
    padding:10px;
}
img
{
    margin:5px;
    border: 1px solid black;
}
</style>
</head>
<body>
<center><h1> #Album title </h1></center>
<p>
EOF

for img in *.jpg;
do
    convert "$img" -resize "100x" *thumbs/$img"
    echo "<a href=\"\$img\" ><img src=\"thumbs/$img\" title=\"\$img\" />
</a>" >> index.html
done

cat <<EOF >> index.html

</p>
</body>
</html>
EOF

echo Album generated to index.html

运行脚本:

$ ./generate_album.sh
Creating album..
Album generated to index.html
```

5.8.3 工作原理

脚本的起始部分用于生成HTML页面的头部。



接下来，脚本将一直到EOF的这部分内容（不包括EOF在内）重定向到index.html；

```
cat <<EOF > index.html
contents...
EOF
```

页面头部包括HTML和样式表单。

for img in *.jpg;对每一个文件进行迭代，并执行相应的操作。

convert "\$img" -resize "100x" "thumbs/\$img"将创建宽度为100px的图像缩略图。

下面的语句会生成所需的标记并将其追加到index.html中：

```
echo "<a href=\"${img}\" ><img src=\"thumbs/${img}\" title=\"${img}\" /></a>" >> index.html
```

最后再用cat添加HTML页脚。

5.8.4 参考

1.5节讲解了EOF和stdin重定向。

5.9 Twitter 命令行客户端

Twitter不仅是最火的微博平台，同时也是最新的在线社会化媒体热点。无论是在Twitter上发帖还是读帖都能让人乐在其中。如果我们要通过命令行来完成这一切呢？其实，写一个基于命令行的Twitter客户端很简单。Twitter有RSS feed，因此可以利用它来实现我们的客户端。来看看怎样做吧。

5.9.1 预备知识

用cURL进行认证，发送Twitter更新并下载RSS feed页面用于解析内容，要做到这一切只需4行代码即可。

5.9.2 实战演练

用curl命令处理Twitter API的Bash脚本如下：

```
#!/bin/bash
#文件名: tweets.sh
#用途: Twitter 客户端基础版

USERNAME="PUT_USERNAME_HERE"
PASSWORD="PUT_PASSWORD_HERE"
COUNT="PUT_NO_OF_TWEETS"

if [[ "$1" != "read" ]] && [[ "$1" != "tweet" ]];
then
```



```

echo -e "Usage: $0 send status_message\n OR\n $0 read\n"
exit -1;
fi

if [[ "$1" = "read" ]];
then
    curl --silent -u $USERNAME:$PASSWORD http://twitter.com/statuses/
    friends_timeline.rss | \
    grep title | \
    tail -n +2 | \
    head -n $COUNT | \
    sed 's:.*<title>\\([^\<)*\\).*:\\n\\1:'

elif [[ "$1" = "tweet" ]];
then
    status=$( echo $@ | tr -d '"' | sed 's:.*tweet //' )
    curl --silent -u $USERNAME:$PASSWORD -d status="$status" http://
    twitter.com/statuses/update.xml > /dev/null
    echo 'Tweeted :)'
fi

```

运行脚本：

```

$ ./tweets.sh tweet Thinking of writing a X version of wall command
"#bash"
Tweeted :)

$ ./tweets.sh read
bot: A tweet line
t3rmln4l: Thinking of writing a X version of wall command #bash

```

5.9.3 工作原理

我们将上面的脚本分为两部分，以便讲解它的运行原理。第一部分是读取tweet消息。从 http://twitter.com/statuses/friends_timeline.rss 下载RSS信息并解析包含 <title> 标记的行，然后用 sed 除去 <title> 和 </title> 标记，这样就得到了所需要的tweet消息内容。利用 head 命令保留 COUNT 条最近的tweet消息，tail -n +2 移除无用的头部文本“Twitter: Timeline of friends”。

第二部分是发送tweet消息，在这一部分中，使用Twitter的API：<http://twitter.com/statuses/update.xml>，通过 curl 的 -d 选项将数据发送到Twitter上。

发送tweet消息时，脚本的 \$1 包含了消息内容。要想获取相关状态信息，我们可以使用 \$@（脚本参数列表），然后从中删除单词“tweet”即可。

5.9.4 参考

- 5.4节讲解了curl命令。
- 3.14节讲解了head与tail。

5.10 基于 Web 后端的定义查询工具

Google 为词汇提供了 Web 定义，可以用搜索查询 “define:WORD” 查找，WORD 即要查找其定义的词汇。我们得用一个图形界面的 Web 浏览器才能检索到这些定义，然而，也可以用脚本自动获取所需要的词汇定义。让我们来看看实现方法。

5.10.1 预备知识

我们用 lynx、sed、awk 和 grep 来编写词汇定义工具。

5.10.2 实战演练

从 Google 搜索检索词汇定义脚本代码如下：

```
#!/bin/bash
#文件名: define.sh
#用途: 一个 Google 定义的前端工具

limit=0
if [ ! $# -ge 1 ];
then
    echo -e "Usage: $0 WORD [-n No_of_definitions]\n"
    exit -1;
fi

if [ "$2" = "-n" ];
then
    limit=$3;
    let limit++
fi

word=$1

lynx -dump http://www.google.co.in/search?q=define:$word | \
awk '/Defini/,/Find defini/' | head -n -1 | sed 's:*\n*::; s:^[ ]*::'
| \
grep -v "[0-9]" | \
awk '{
    if ( substr($0,1,1) == "*" )
    { sub("++count*."); };
    print
}' > /tmp/$$.txt

echo
if [ $limit -ge 1 ];
then

cat /tmp/$$.txt | sed -n "/^1$/, /${limit}/p" | head -n -1

else
```

```
cat /tmp/$$.txt;
```

```
fi
```

运行脚本：

```
$ ./define.sh hack -n 2
1. chop: cut with a hacking tool
2. one who works hard at boring tasks
```

5.10.3 工作原理

让我们来研究一下定义解析器的核心部分。Lynx用来获取纯文本形式的网页。定义页面的URL是[http://www.google.co.in/search?q=define:\\$word](http://www.google.co.in/search?q=define:$word)。然后我们减少“Definitions on web”与“Find definitions”之间的文本。所有的定义都出现在这些文本行中(`awk '/Defini/,/Find defini/'`)。

`'s:*\n*:'` 用来将 * 替换成 `\n*`，这是为了在每个定义之间插入一个换行符。`s:^()*::` 用来移除行首多余的空格。超链接在 lynx 输出中被标记成 `[数字]` 的形式。这些行可以用 `grep -v` (输出不匹配的行) 移除。然后用 `awk` 将行首的 * 替换成数字，为每一个定义分配一个序号。如果我们在脚本中使用了 `-n count` (`count` 在这里表示希望打印的定义数)，那就只输出 `count` 项定义。因此，用 `awk` 打印从编号1到编号 `count` 的定义 (因为我们已经将 * 替换成了序列号，所以这一步实现起来就容易多了)。

5.10.4 参考

- 4.6节讲解了 `sed` 命令。
- 4.7节讲解了 `awk` 命令。
- 4.3节讲解了 `grep` 命令。
- 5.3节讲解了 `lynx` 命令。

5.11 查找网站中的无效链接

我曾目睹过人们纯手动检查网站上的每一个页面以便查找出无效的链接。这对于一些页面不多的小网站来说还有可能。如果页面数量增多，纯手动的话就不现实了。要是能够将查找无效链接的工作自动化，那可轻松太多了。而我们也确实可以借助HTTP处理工具来查找无效的链接。来看看这是怎样实现的。

5.11.1 预备知识

要识别链接并从中找出无效链接，我们可以使用 `lynx` 和 `curl`。`lynx` 有一个选项 `-traversal`，能够以递归方式访问网站页面并建立网站中所有超链接的列表。我们可以用 `curl` 验证每一个链接的有效性。

5.11.2 实战演练

用curl查找网页上无效链接的Bash脚本如下：

```
#!/bin/bash
#文件名: find_broken.sh
#用途: 查找网站中的无效链接

if [ $# -eq 2 ];
then
    echo -e "$Usage $0 URL\n"
    exit -1;
fi

echo Broken links:

mkdir /tmp/${$.lynx}

cd /tmp/${$.lynx}

lynx -traversal $1 > /dev/null
count=0;
sort -u reject.dat > links.txt

while read link;
do
    output=`curl -I $link -s | grep "HTTP/*OK*"`;
    if [[ -z $output ]];
    then
        echo $link;
        let count++
    fi
done < links.txt

[ $count -eq 0 ] && echo No broken links found.
```

5

5.11.3 工作原理

lynx -traversal URL会在工作目录下生成数个文件。这其中包括reject.dat，该文件包含网站中所有的链接。sort -u用来建立一个不包含重复项的列表。然后我们迭代每一个链接，并通过curl -I检验接收到的响应头部。如果响应头部的第一行包含HTTP/1.0 200 OK，就表示该链接正常。显示其他响应信息则表示该链接失效，并将其输出到stdout。

5.11.4 参考

- 5.3节讲解了lynx命令。
- 5.4节讲解了curl命令。

5.12 跟踪网站变更

对于Web开发人员和用户来说，能够跟踪网站变更不无益处。但定期手动检查既麻烦也不现实，因此，我们可以编写一个定期运行的变更跟踪器（change tracker）。一旦发生变更，跟踪器便会发出声音或发送提示信息。下面看一看如何编写一个基础的网站变更跟踪器。

5.12.1 预备知识

用Bash脚本跟踪网站变更意味着要在不同的时间检索网站，然后用diff命令进行比对。我们可以使用curl和diff来实现这些步骤。

5.12.2 实战演练

结合各种命令来跟踪网页变更的Bash脚本如下：

```
#!/bin/bash
#文件名: change_track.sh
#用途: 跟踪网页更新

if [ $# -eq 2 ];
then
    echo -e "$Usage $0 URL\n"
    exit -1;
fi

first_time=0
# 非首次运行

if [ ! -e "last.html" ];
then
    first_time=1
    # 首次运行
fi

curl --silent $1 -o recent.html

if [ $first_time -ne 1 ];
then
    changes=$(diff -u last.html recent.html)
    if [ -n "$changes" ];
    then
        echo -e "Changes:\n"
        echo "$changes"
    else
        echo -e "\nWebsite has no changes"
    fi
else
    echo "[First run] Archiving..."
```



```
fi
```

```
cp recent.html last.html
```

分别观察一下网页发生更新与没有更新时脚本track_changes.sh的输出。

□ 第一次运行

```
$ ./track_changes.sh http://web.sarathlakshman.info/test.html
[First run] Archiving..
```

□ 第二次运行

```
$ ./track_changes.sh http://web.sarathlakshman.info/test.html
Website has no changes
```

□ 在网页更新后，第三次运行

```
$ ./test.sh http://web.sarathlakshman.info/test_change/test.html
Changes:

--- last.html 2010-08-01 07:29:15.000000000 +0200
+++ recent.html 2010-08-01 07:29:43.000000000 +0200
@@ -1,3 +1,4 @@
<html>
+added line :)
<p>data</p>
</html>
```

5

5.12.3 工作原理

脚本用 [! -e "last.html"]; 检查脚本本身是否是首次运行。如果last.html不存在，那就意味着这是首次运行，因此要下载网页并将其复制为last.html。

如果不是第一次运行，那么脚本应该下载一个新的网页副本 (recent.html)，然后用diff检查差异。如果有变化，则打印出变更信息并将recent.html复制成last.html。

5.12.4 参考

5.4节讲解了curl命令。

5.13 以 POST 方式发送网页并读取响应

POST和GET是HTTP协议中用于发送或检索信息的两种请求类型。在GET请求方式中，我们利用网页的URL来发送参数（“名称-值”对）。而在POST请求方式中，我们就不这么做了。POST方式用于提交表单，诸如提交用户名、密码以及检索登录页面等。

当编写基于网页检索的脚本时，用POST方式发送页面相当常见。让我们来看一看POST的使用方法。通过POST方式发送数据并检索输出自动化HTTP GET和POST请求，在编写解析网站数据的shell脚本的过程中，是一项非常重要的任务。

5.13.1 预备知识

cURL和Iwget都可以通过各种参数来处理POST请求。请求内容能够以“名称-值”对的形式传递给命令。

5.13.2 实战演练

来看看如何用curl发送POST请求并读取网站的HTML响应:

```
$ curl URL -d "postvar=postdata2&postvar2=postdata2"
```

在一个用于提交诸如主机名和用户名等当前用户信息的网站中 (http://book.sarathlakshman.com/lsc/mlogs/), 假设其主页有HOSTNAME和USER两个字段和一个SUBMIT按钮, 当用户输入主机名和用户名, 点击SUBMIT按钮, 输入的信息就会储存在网站中。这个过程可以通过一行curl命令将POST请求提交自动化来实现。如果你查看网站的源代码 (使用网页浏览器的查看源代码选项), 你会发现一个与下面类似的HTML表单:

```
<form action="http://book.sarathlakshman.com/lsc/mlogs/submit.php"
method="post" >

<input type="text" name="host" value="HOSTNAME" >
<input type="text" name="user" value="USER" >
<input type="submit" >

</form>
```

其中 http://book.sarathlakshman.com/lsc/mlogs/submit.php 是目标URL。当用户输入信息并点击submit按钮, 名为host和user的输入内容就以POST请求的方式被发送到submit.php中, 然后对应的响应页面被返回到浏览器。

我们可以按照下面的方法自动化POST请求:

```
$ curl http://book.sarathlakshman.com/lsc/mlogs/submit.php -d "host=test-
host&user=slynux"
<html>
You have entered :
<p>HOST : test-host</p>
<p>USER : slynux</p>
<html>
```

curl返回响应页面。

-d选项用于发送请求。-d的字符串参数和GET请求在语义上类似。每对var=value之间用&分隔。



-d的参数内容应该以引用的形式给出。如果不使用引用, & 会被shell解读为该命令应该作为后台进程运行。

5.13.3 补充内容

让我们再看看如何用cURL和wget执行POST请求。

1. 在curl中使用POST

可以在curl中用 -d 或 --data 以POST方式发送数据：

```
$ curl --data "name=value" URL -o output.html
```

如果需要发送多个变量，用 & 分隔变量。需要注意的是：如果使用了 &，那么该“名称-值”对应以引用的形式出现，否则shell会将其看做是用于后台进程的特殊符号。例如：

```
$ curl -d "name1=val1&name2=val2" URL -o output.html
```

2. 用wget以POST方式发送数据

用wget的--post-data "string" 可以POST方式发送数据。例如：

```
$ wget URL --post-data "name=value" -O output.html
```

其中的“名称-值”对要采用和cURL同样的格式。

5.13.4 参考

- 5.4节讲解了curl命令。
- 5.2节讲解了wget命令。



本章内容

- ☐ 用tar归档
- ☐ 用cpio归档
- ☐ 用gunzip或gzip压缩
- ☐ 用bunzip或bzip压缩
- ☐ 用lzma压缩
- ☐ 用zip归档和压缩
- ☐ 超高压缩率的squashfs文件系统
- ☐ 用标准算法加密文件和文件夹
- ☐ 用raync备份系统快照
- ☐ 用git备份版本控制
- ☐ 用dd克隆磁盘

6.1 简介

提取快照和备份数据都是我们的日常工作。就服务器或大型数据存储系统而言，定期备份不可小视。我们可以通过shell脚本来实现备份自动化。归档和压缩对于系统管理员或普通用户来说平日里都少不了要用到。有许多不同的压缩格式，要结合不同的使用方法才可以获得最佳的压缩效果。加密是另一种保护数据的常用方法。为了减少加密数据的大小，文件在加密之前通常都要先进行归档和压缩。有很多标准加密算法可以使用，并且它们都有相应的shell工具。本章考查了这方面的各类攻略，包括创建和维护文件或文件夹归档、压缩格式以及加密技术。接下来让我们看看这些内容。

6.2 用tar 归档

tar命令可以用来归档文件。不过，它最初是设计用来将数据存储于磁带上（tape archives, tar）。可以用tar将多个文件和文件夹保存为单个文件，同时还能保留所有的文件属性，如所有者、权限等。由该命令创建的文件通常称为tarball。

6.2.1 预备知识

所有类UNIX操作系统都默认包含tar命令。它的语法简单，并且文件格式具备可移植性。让我们来看看命令的用法。

tar 包含一个参数列表: A、c、d、r、u、x、f 和 v。其中每一个参数都可以依据不同的用途单独使用。

6.2.2 实战演练

按照下面的语法, 用 tar 对文件进行归档:

```
$ tar -cf output.tar [SOURCES]
```

例如:

```
$ tar -cf output.tar file1 file2 file3 folder1 ..
```

命令中的 -c 代表“创建文件”(creat file), -f 代表“指定文件名”(specify filename)。

我们可以使用文件名列表或者诸如 *.txt 这类通配符来指定文件夹和文件名作为 SOURCES。

上面的命令将原文件归档为单个文件 output.tar。

文件名必须紧跟在 -f 之后, 并且应该是参数组中的最后一项 (例如, -cvvf filename.tar 和 -tvvf filename.tar)。

我们不能传递上百个文件或文件夹作为 tar 的参数, 毕竟参数数量不是无限制的。如果有很多文件需要归档的话, 使用追加 (append) 选项要更安全些。

6.2.3 补充知识

让我们再来看看 tar 命令的其他特性。

1. 向归档文件中添加文件

有时候, 我们可能需要向已存在的归档文件再添加一些文件 (其用例情境包括: 当有成百上千个文件需要归档且这些文件又无法作为命令行参数在一行中指定时)。

追加选项: -r

要向已存在的归档文件添加一个文件, 可以使用:

```
$ tar -rvf original.tar new_file
```

用下面的方法列出归档文件中的内容:

```
$ tar -tf archive.tar
yy/lib64/
yy/lib64/libfakeroot/
yy/sbin/
```

如果需要在归档或列出归档内容的过程中获知更多细节, 可以使用 -v 或 -vv 选项。它们统称为 verbose (v), 它们允许在终端中输出更详细的信息, 借助 verbose, 你能够打印诸如文件权限、所有者所属组、文件修改日期等细节。

例如:

```
$ tar -tvvf archive.tar
drwxr-xr-x slynux/slynux 0 2010-08-06 09:31 yy/
```

```
drwxr-xr-x slynux/slynux 0 2010-08-06 09:39 yy/user/
drwxr-xr-x slynux/slynux 0 2010-08-06 09:31 yy/user/lib64/
```

2. 从归档文件中提取文件或文件夹

下面的命令将归档文件的内容提取到当前目录中：

```
$ tar -xf archive.tar
```

选项-x表示提取(exact)。

使用-x时，tar命令将归档文件中的内容提取到当前目录。我们也可以使用选项-C指定需要将文件提取到哪个目录：

```
$ tar -xf archive.tar -C /path/to/extraction_directory
```

这个命令将归档文件的内容提取到指定目录中。它提取的是归档文件中的全部内容。我们可以通过将文件名指定为命令行参数来提取特定的文件：

```
$ tar -xvf file.tar file1 file4
```

上面的命令只提取file1和file4，同时忽略掉其他文件。

3. 在tar中使用stdin和stdout

进行归档时，我们可以将stdout指定为输出文件，这样另一个命令就可以通过管道将它作为stdin，并进行其他处理或提取内容。

当通过安全shell(Secure Shell, SSH)连接传输数据时，这招很管用。例如：

```
$ mkdir ~/destination
$ tar -cf - file1 file2 file3 | tar -xvf - -C ~/destination
```

在上面的例子中，首先对file1、file2和file3进行归档，然后将它们提取到~/destination中。在这个命令中：

□ -f指定stdout作为归档文件(当使用-c选项时)

□ -f指定stdin用于提取内容(当使用-x选项时)

4. 拼接两个归档文件

我们可以用-A选项轻松地合并多个tar文件。

假设我们现在有两个tar文件：file1.tar和file2.tar。我们可以按照下面的方法将file2.tar的内容合并到file1.tar中：

```
$ tar -Af file1.tar file2.tar
```

查看内容，验证操作是否成功：

```
$ tar -tvf file1.tar
```

5. 通过检查时间戳来更新归档文件中的内容

添加选项可以将指定的任意文件加入到归档文件中。如果同名文件已经存在，那么结果就是在归档文件中包含了两个同名的文件。我们可以用更新选项-u指明：只有比归档文件中的同名文件更新的时候才进行添加。

```
$ tar -tf archive.tar
filea
fileb
filec
```

上面的命令将列出归档文件中的内容。

为了做到只有在filea的文件内容修改时间更新 (newer) 的时候才对它进行添加, 可以使用:

```
$ tar -uvvf archive.tar filea
```

如果我们想要追加的filea和归档文件中的filea有相同的时间戳, 那么不执行任何操作。这时可用touch命令修改文件的时间戳, 然后再用tar命令:

```
$ tar -uvvf archive.tar filea
-rw-r--r-- slynux/slynux      0 2010-08-14 17:53 filea
```

因为时间戳比归档文件中的同名文件更新, 因此执行追加操作。

6. 比较归档文件与文件系统中的内容

有时候, 如果能够知道归档文件中的文件与文件系统中的同名文件是否有差别, 对于我们而言还是有帮助的。选项 -d 可以打印出两者之间的差别:

```
$ tar -df archive.tar filename1 filename2 ...
```

例如:

```
$ tar -df archive.tar afile bfile
afile: Mod time differs
afile: Size differs
```

7. 从归档文件中删除文件

我们可以用 --delete 选项从给定的归档文件中删除文件。例如:

```
$ tar -f archive.tar --delete file1 file2 ..
```

来看另外一个例子:

```
$ tar -tf archive.tar
filea
fileb
filec
```

我们也可以使用下面的语法:

```
$ tar --delete --file archive.tar [FILE LIST]
```

例如:

```
$ tar --delete --file archive.tar filea
$ tar -tf archive.tar
fileb
filec
```

8. 压缩tar归档文件

tar 命令只能用来对文件进行归档, 它并不具备压缩功能。出于这个原因, 多数用户在使用归档文件时都会对文件采用某种形式的压缩, 这样就能够大量降低文件的大小。归档文件通常被

压缩成下列格式之一：

- ☐ file.tar.gz
- ☐ file.tar.bz2
- ☐ file.tar.lzma
- ☐ file.tar.lzo

不同的tar选项可以用来指定不同的压缩格式。

- ☐ -j指定bunzip2格式；
- ☐ -z指定gzip格式；
- ☐ --lzma指定lzma格式。

这些格式会在随后专门讲解压缩的攻略中讨论。

可以不明确指定上面那些特定的选项来使用压缩功能。tar可以通过查看输出或输入文件的扩展名来进行压缩。为了让tar支持根据扩展名自动进行压缩，使用-a或--auto-compress选项。

9. 从归档中排除部分文件

通过指定样式可以从归档中排除部分文件。用--exclude [PATTERN] 排除匹配通配符样式的文件。

例如，排除所有的.txt文件：

```
$ tar -cf arch.tar * --exclude "*.txt"
```



样式应该使用双引号来引用。

也可以将需要排除的文件列表放入文件中，同时配合选项 -X：

```
$ cat list
filea
fileb

$ tar -cf arch.tar * -X list
```

这样就将filea和fileb从归档中排除了。

10. 排除版本控制目录

我们通常使用tar归档文件来分发源代码。大多数源代码都是使用版本控制系统进行维护，如subversion、git、mercurial、cvs等。版本控制系统中的代码目录包含用来管理版本的特殊目录，如svn或.git。但源代码本身并不需要这些目录，所以不应该将它们包含在源代码的tar归档文件中。

为了在归档时排除版本控制相关的文件和目录，可以使用tar的--exclude-vcs选项。例如：

```
$ tar --exclude-vcs -czvzf source_code.tar.gz eye_of_gnome_svn
```

11. 打印总字节数

有时候我们需要打印出归档了多少字节。用 -totals就可以在归档完成之后打印出总归档

字节数：

```
$ tar -cf arc.tar * --exclude "*.txt" --totals
Total bytes written: 20480 (20KiB, 12MiB/s)
```

6.2.4 参考

- 6.4节讲解了gzip命令。
- 6.5节讲解了bzip2命令。
- 6.6节讲解了lzma命令。

6.3 用 cpio 归档

cpio是类似于tar的另一种归档格式。它用来将多个文件和文件夹存储为单个文件，同时还能保留所有的文件属性，如权限、文件所有权等。不过，cpio用得不如tar那么多。它多用于RPM软件包、Linux内核的initramfs文件等。本攻略将给出几个cpio的用例。

实战演练

cpio通过stdin获取输入文件名，并将归档文件写入stdout。我们必须将stdout重定向到一个文件，以接收cpio的输出。

创建测试文件：

```
$ touch file1 file2 file3
```

将测试文件按照下面的方法进行归档：

```
$ echo file1 file2 file3 | cpio -ov > archive.cpio
```

在上面的命令中：

- -o指定了输出；
- -v用来打印归档文件列表。

6



通过cpio，我们能够利用文件的绝对路径进行归档。/usr/somedir就是一个绝对路径，因为它以根目录(/)作为路径的起始。相对路径不以/开头，而是以当前目录作为起始点。例如，test/file表示有一个目录test，而file位于test目录中。

当进行提取时，cpio会使用绝对路径提取内容。但是tar会移去绝对路径开头的/，将之转换为相对路径。

要列出cpio归档文件中的内容，使用下列命令：

```
$ cpio -it < archive.cpio
```

这个命令会列出给定cpio归档文件中的所有内容。它从stdin中读取文件。在该命令中：

- `-i`用于指定输入;
 - `-t`表示列出归档文件中的内容。
- 要从cpio归档文件中提取文件, 可以使用:

```
$ cpio -id < archive.cpio
```

这里的`-d`用来表示提取内容。

cpio在覆盖文件时不会发出提示。如果绝对路径上的文件已经存在, cpio会将其替换。它不会像tar那样将文件提取到当前目录中。

6.4 用 gunzip 或 gzip 压缩

gzip是GNU/Linux平台下常用的压缩格式。gzip、gunzip、zcat都可以处理gzip压缩文件类型。gzip只能够压缩单个文件, 而无法对目录和多个文件进行归档, 因此我们将这个任务交给tar, 然后再用gzip进行压缩。如果指定多个文件作为输入, 将会生成多个单独的压缩文件(.gz)。让我们来看看gzip的使用方法。

6.4.1 实战演练

要使用gzip压缩文件, 可以使用下面的命令:

```
$ gzip filename
$ ls
filename.gz
```

gzip会删除原文件并生成一个压缩文件filename.gz。

将gzip文件解压缩的方法如下:

```
$ gunzip filename.gz
```

gunzip会删除filename.gz并生成filename.gz的未压缩形式。

要列出压缩文件的属性信息, 可以使用:

```
$ gzip -l test.txt.gz
compressed  uncompressed  ratio  uncompressed_name
35          6              -33.3% test.txt
```

gzip命令可以从stdin中读入文件, 也可以将压缩文件写出到stdout。

从stdin读入并写出到stdout:

```
$ cat file | gzip -c > file.gz
```

选项`-c`用来将输出指定到stdout。

我们可以指定gzip的压缩级别。用`--fast`或`--best`选项分别提供最低或最高的压缩比。

6.4.2 补充内容

gzip命令通常与其他命令结合使用。它还有一些高级选项可以用来指定压缩比。让我们来

看看gzip的这些特性。

1. 压缩归档文件

我们通常将gzip与归档文件结合使用。当进行归档及提取的时候，可以用tar命令的 -z选项来压缩归档文件。

你可以按照下面的方法创建经由gzip压缩过的归档文件：

□ 方法 1

```
$ tar -czvf archive.tar.gz [FILES]
```

或者

```
$ tar -cavvf archive.tar.gz [FILES]
```

选项 -a指定从文件扩展名自动判断压缩格式。

□ 方法 2

首先，创建一个归档文件：

```
$ tar -cvvf archive.tar [FILES]
```

压缩归档文件：

```
$ gzip archive.tar
```

如果有多个文件（上千个）需要归档及压缩，我们可以采用方法2，只需在此基础上稍作变动。将多个文件作为命令行参数传递给tar的问题在于：tar只能从命令行中接受有限个文件。要解决这个问题，我们可以在一个循环中使用添加选项（-r）来逐个添加文件：

```
FILE_LIST="file1 file2 file3 file4 file5"
for f in $FILE_LIST;
do
tar -rvf archive.tar $f
done
gzip archive.tar
```

如果要提取经由gzip压缩的归档文件中的内容，可以使用：

□ -x用于提取内容；

□ -z表示采用gzip格式。

或者

```
$ tar -xavvf archive.tar.gz -C extract_directory
```

在上面的命令中，选项 -a用于自动检测压缩格式。

2. zcat——无需解压缩，直接读取gzip格式文件

zcat命令无需解压缩，直接就可以从 .gz文件中提取内容，并打印到stdout。按照下面的方法，在保持 .gz文件不变的情况下，将从中提取出的文件打印出来：

```
$ ls
test.gz

$ zcat test.gz
```

```
A test file
# 文件test包含了一行文本"A test file"
```

```
$ ls
test.gz
```

3. 压缩率

我们可以指定压缩率，压缩率有9级，其中：

- 1级的压缩率最低，但是压缩速度最快；
- 9级的压缩率最高，但是压缩速度最慢。

你可以按照下面的方法指定压缩比：

```
$ gzip -9 test.img
```

它将按照最大的压缩比对文件进行压缩。

6.4.3 参考

6.2节讲解了tar命令。

6.5 用 bzip2 或 bzip 压缩

bzip2是另一种与gzip类似的压缩技术。bzip2通常能够生成比gzip更小（压缩比更高）的文件。所有Linux发行版都包含了这个工具。让我们看看bzip2的用法。

6.5.1 实战演练

用bzip2进行压缩：

```
$ bzip2 filename
$ ls
filename.bz2
```

bzip2会删除原文件并生成名为filename.bz2的压缩文件。

解压bzip2格式的文件：

```
$ bunzip2 filename.bz2
```

bzip2会删除filename.bz2并生成filename的未压缩形式。

bzip2可以从stdin中读入文件，也可以将压缩文件写出到stdout。

要从stdin读入并写出到stdout，使用：

```
$ cat file | bzip2 -c > file.tar.bz2
```

-c用于将输出指定到stdout。

我们通常将bzip2与归档文件结合使用。当进行归档及提取的时候，可以用tar命令的-j选项来压缩归档文件。

下面的方法可以创建经由bzip2压缩过的归档文件。

❑ 方法 1

```
$ tar -cjvfv archive.tar.bz2 [FILES]
```

或者

```
$ tar -cavfv archive.tar.bz2 [FILES]
```

选项-a指定由文件扩展名自动判断压缩格式。

❑ 方法 2

首先创建一个归档文件：

```
$ tar -cvfv archive.tar [FILES]
```

压缩归档文件：

```
$ bzip2 archive.tar
```

如果我们要对成百上千个文件进行归档，那上面的命令可能就不管用了。要解决这个问题，需要使用-r选项并借助一个循环来将文件逐个添加到归档文件中。请参考6.4节中类似的解决方法。

提取bzip2压缩格式的归档文件中的内容，可以使用：

```
$ tar -xjvfv archive.tar.bz2 -C extract_directory
```

其中：

❑ -x用于提取内容；

❑ -j表示采用bzip2格式；

❑ -C指定用于存放提取文件的目录。

或者，使用下面的命令：

```
$ tar -xavfv archive.tar.bz2 -C extract_directory
```

-a会自动检测压缩格式。

6

6.5.2 补充内容

bunzip还有其他选项，可用以执行多种功能。让我们简单了解几个。

1. 保留输入文件

使用bzip2或bunzip2时，它会删除输入文件并生成压缩过的输出文件。我们可以使用选项-k来避免删除输入文件。例如：

```
$ bunzip2 test.bz2 -k
$ ls
test test.bz2
```

2. 压缩率

我们可以指定压缩率，从1级到9级（其中，第1级压缩率最低，但是压缩速度最快；第9级压缩率最高，但是速度要慢得多）。

例如：

```
$ bzip2 -9 test.img
```

该命令规定的是最大压缩率。

6.5.3 参考

6.2节讲解了tar命令。

6.6 用 lzma 压缩

与gzip或bzip2相比，lzma是一个相对较新的压缩工具。它提供了比gzip或bzip2更好的压缩率。由于在大多数Linux发行版上都没有预装lzma，你得使用软件包管理器自行安装。

6.6.1 实战演练

使用lzma进行压缩：

```
$ lzma filename
$ ls
filename.lzma
```

lzma会删除原文件并生成名为filename.lzma的压缩文件。

解压缩lzma文件：

```
$ unlzma filename.lzma
```

该命令会删除filename.lzma并生成该文件的未压缩形式。

lzma命令可以从stdin中读入文件，也可以将压缩文件写出到stdout。

从stdin读入并写出到stdout：

```
$ cat file | lzma -c > file.lzma
```

-c用于将输出指定到stdout。

我们通常将lzma与归档文件结合使用。当进行归档及提取的时候，可以使用tar命令的-lzma选项来压缩归档文件。

有两种方法可以创建lzma归档文件。

□ 方法 1

```
$ tar -cvvf --lzma archive.tar.lzma [FILES]
```

或者

```
$ tar -cavvf archive.tar.lzma [FILES]
```

选项-a指定由文件扩展名自动判断压缩格式。

□ 方法 2

首先创建一个归档文件：

```
$ tar -cvvf archive.tar [FILES]
```

压缩归档文件：

```
$ lzma archive.tar
```

如果我们要对成百上千个文件进行归档，上面的命令可能不太管用。要解决这个问题，需要使用-r选项并借助一个循环来将文件逐个追加到归档文件中。请参考6.4节中类似的解决方法。

6.6.2 补充内容

来看看lzma的其他选项。

1. 提取lzma归档文件中的内容

要将使用lzma压缩过的归档文件中的内容提取到指定的目录，可以使用：

```
$ tar -xvzf --lzma archive.tar.lzma -C extract_directory
```

其中，-x用于提取内容，--lzma指定使用lzma对归档文件进行解压缩。

我们也可以使用：

```
$ tar -xavzf archive.tar.lzma -C extract_directory
```

选项-a指定由文件扩展名自动判断压缩格式。

2. 保留输入文件

使用lzma或unlzma时，它们会删除输入文件并生成压缩过的输出文件。不过我们可以用选项-k来避免删除输入文件。例如：

```
$ lzma test.bs2 -k
$ ls
test.bs2.lzma
```

3. 压缩率

我们可以指定压缩率，从1级到9级（其中，第1级压缩率最低，但是压缩速度最快，第9级压缩率最高，但是速度要慢得多）。

你可以按照下面的方法指定压缩率：

```
$ lzma -9 test.img
```

该命令规定的是最大压缩率。

6.6.3 参考

6.2节讲解了tar命令。

6.7 用 zip 归档和压缩

ZIP作为一种流行的压缩格式，在很多平台中都可以看到它的身影。在Linux下，它的应用不如gzip或bzip2那么广泛，但是Internet上的文件通常都采用这种格式。

实战演练

对归档文件采用ZIP格式进行压缩：

```
$ zip archive_name.zip [SOURCE FILES/DIRS]
```

例如：

```
$ zip file.zip file
```

该命令会生成file.zip。

对目录和文件进行递归操作：

```
$ zip -r archive.zip folder1 file2
```

其中，-r用于指定递归指定操作。

和lzip、gzip以及bzip2不同的是，zip在完成压缩之后不会删除原文件。在这方面，它和tar类似，而除此之外，zip还拥有tar所不具备的压缩功能。

要从ZIP文件中提取内容，可以使用：

```
$ unzip file.zip
```

在完成提取操作之后，unzip并不会删除file.zip（这一点与unlzip和gunzip不同）。

如果需要更新归档文件中的内容，使用选项 -u：

```
$ zip file.zip -u newfile
```

从压缩归档文件中删除内容，则使用-d：

```
$ zip -d arc.zip file.txt
```

列出归档文件中的内容：

```
$ unzip -l archive.zip
```

6.8 超高压压缩率的 squashfs 文件系统

squashfs是一种只读型的超高压压缩率文件系统，这种文件系统能够将2GB-3GB的数据压缩成一个700MB的文件。你有没有想过Linux Live CD是怎样运行的？当Live CD启动后，它会加载一个完整的Linux环境。这就是利用了一种被称为squashfs的只读型压缩文件系统。它将根文件系统保存在一个压缩过的文件系统文件中。这个文件可以使用环回的形式来挂载并对其中的文件进行访问。因此当进程需要某些文件，可以将它们解压，然后载入内存中使用。如果需要构建一个定制的Live OS，或是需要使用超高压压缩率的文件并且无需解压就可以访问文件，那么squashfs的相关知识就能派得上用场。要解压个头较大的压缩文件，那可得花上一阵工夫。但如果将文件以环回形式挂载，速度就会飞快，因为只有出现访问请求的时候，对应的那部分压缩文件才会被解压缩。而普通的解压缩方式是首先解压缩所有的数据。让我们来看看如何使用squashfs。

6.8.1 预备知识

如果你有一张Ubuntu CD, 可以在CDRom ROOT/casper/filesystem.squashfs中找到文件.squashfs。squashfs在内部采用了gzip和lzma这类压缩算法。所有最新的Linux发行版都支持squashfs。但要想创建squashfs文件, 需要额外安装squashfs-tools软件包。

6.8.2 实战演练

通过添加源目录和文件, 创建一个squashfs文件:

```
$ mksquashfs SOURCES compressedfs.squashfs
```

SOURCES部分可以使用通配符或文件、目录路径。

例如:

```
$ sudo mksquashfs /etc test.squashfs
Parallel mksquashfs: Using 2 processors
Creating 4.0 filesystem on test.squashfs, block size 131072.
[=====] 1867/1867 100%
```

还有更多的细节信息会出现在终端上。由于版面的限制, 这里就不再列出这些信息了

要挂载squashfs文件, 利用环回形式进行挂载:

```
# mkdir /mnt/squash
# mount -o loop compressedfs.squashfs /mnt/squash
```

你可以通过访问/mnt/squashfs复制其中的内容。

6.8.3 补充内容

可以指定额外的参数来创建squashfs文件系统。让我们来看一些其他的命令选项。

在创建squashfs文件时排除部分文件

创建squashfs文件时, 我们可以排除部分文件, 这些文件可以用文件列表的形式指定, 也可以用通配符来指定。

使用选项 -e, 将需要排除的文件列表以命令行参数的方式来指定。例如:

```
$ sudo mksquashfs /etc test.squashfs -e /etc/passwd /etc/shadow
```

其中, 选项 -e用于排除文件passwd和shadow。

也可以将需要排除的文件名写入文件, 然后用 -ef指定该文件:

```
$ cat excludelist
/etc/passwd
/etc/shadow
```

```
$ sudo mksquashfs /etc test.squashfs -ef excludelist
```

如果希望在排除文件列表中使用通配符, 那么可以使用 -wildcard选项。

6.9 加密工具与散列

加密技术主要用于防止数据遭受未经授权的访问。加密算法有很多，我们会着重讲解那些常见的标准加密算法。在Linux环境下有些工具可以用来执行加密和解密。有时我们使用加密算法散列值来验证数据的完整性。本节将介绍一些常用的加密工具以及这些工具所涉及的算法。

实战演练

让我们看看crypt、gpg、base64、md5sum、shasum以及openssl的用法。

□ crypt

crypt是一个简单的加密工具，它从stdin接受一个文件以及口令（passphrase）作为输入，然后将加密数据输出到output。

```
$ crypt <input_file> output_file
Enter passphrase:
```

它会要求输入一个口令。而我們也可以通过命令行参数来提供口令。

```
$ crypt PASSPHRASE < input_file > encrypted_file
```

如果需要解密文件，可以使用：

```
$ crypt PASSPHRASE -d < encrypted_file > output_file
```

□ gpg

gpg (GNU privacy guard, GNU 隐私保护) 是一种应用广泛的加密方案，它采用密钥签名技术保护文件内容，只有经过认证的用户才能访问数据。我们对gpg签名早已耳熟能详。不过，gpg的技术细节已经超出了本书讨论的范围，在这里我们只学习如何用它对文件进行加密和解密。

用gpg加密文件：

```
$ gpg -c filename
```

这条命令采用交互方式读取口令，并生成filename.gpg。

解密gpg文件：

```
$ gpg filename.gpg
```

这条命令读取口令，然后对文件进行解密。

□ Base64

Base64是一组类似的编码方案（encoding scheme），它通过将ASCII字符转换成以64为基数的形式（radix-64 representation）来用ASCII字符串描述二进制数据。base64可以用来对编码和解码Base64字符串。

要将文件编码为Base64格式，可以使用：

```
$ base64 filename > outputfile
```

或者

```
$ cat file | base64 > outputfile
```

base64可以从stdin中进行读取。

按照下面的方法解码Base64数据:

```
$ base64 -d file > outputfile
```

或者

```
$ cat base64_file | base64 -d > outputfile
```

□ md5sum与sha1sum

md5sum与sha1sum都是单向散列算法(unidirectional hash algorithm), 均无法逆推出原始数据。它们通常用于验证数据完整性或为特定数据生成唯一的密钥, 因为通过分析文件内容, 它们可以为每个文件生成一个唯一的密钥。

```
$ md5sum file
8503063d5488c3080d4800ff50850dc9 file
```

```
$ sha1sum file
1ba02b66e2e557fede8f61b7df282cd0a27b816b file
```

这种类型的散列算法是存储密码的理想方案。密码使用其对应的散列值来存储。如果某个用户需要进行认证, 读取该用户提供的密码并转换成散列值, 然后将其与之前存储的散列值进行比对。如果相同, 用户就通过认证, 被允许访问; 否则, 就会被拒绝访问。将密码以明文形式存储是件非常冒险的事, 它面临密码泄漏的风险。

□ shadowlike散列(salted散列)

让我们看看如何为密码生成shadowlike salted散列。

在Linux中, 用户密码是以散列值形式存储在文件/etc/shadow中的。该文件中一行典型的内容类似于这样:

```
test:$6$fG4eWdUi$ohTKOLeUzNk77.4S8MrYe07NTRV4M3LrJnZP9p.qc1bR5c.EcOruzPXfEululo
BFUa18ENRH7F70zhodas3cR.:14790:0:99999:7:::
```

在这行中, \$6\$fG4eWdUi\$ohTKOLeUzNk77.4S8MrYe07NTRV4M3LrJnZP9p.qc1bR5c.EcOruzPXfEululoBFUa18ENRH7F70zhodas3cR是密码对应的shadow散列值。

某些情况下, 我们可能需要编写一些重要的管理脚本, 这些脚本也许需要编辑密码, 或是需要用shell脚本手动添加用户。这时我们必须生成shadow密码字符串, 并向shadow文件中写入类似于上面的文本行。下面让我们看看如何用openssl生成shadow密码。

shadow密码通常都是salted密码。所谓SALT就是额外的一个字符串, 用来起一个混淆的作用, 是加密更加不易被破解。salt由一些随机位组成, 被用作密钥生成函数的输入之一, 以生成密码的salted散列值。

关于 salt 的更多细节信息，请参考维基百科页面 [http://en.wikipedia.org/wiki/Salt_\(cryptography\)](http://en.wikipedia.org/wiki/Salt_(cryptography))。

```
$ openssl passwd -1 -salt SALT_STRING PASSWORD
$1$SALT_STRING$323VkwkSLRuhbt1zk8eUG.
```

它将 SALT_STRING 替换为随机字符串，同时也将 PASSWORD 替换成你想要使用的密码。

6.10 用 rsync 备份系统快照

备份数据是多数系统管理员的日常工作。我们可能要备份 Web 服务器上的数据，也可能要备份远端数据。rsync 可以对位于不同位置的文件和目录进行备份，它可以借助差异计算以及压缩技术来最小化数据传输量。相对于 cp 命令，rsync 的优势在于前者使用了高效的差异算法 (difference algorithm)。另外，它还支持网络数据传输。在进行复制的同时，rsync 会比较源端和目的端的文件，只有当文件有更新时才进行复制。rsync 也支持压缩、加密等多种特性。让我们看看 rsync 的使用方法。

6.10.1 实战演练

将源目录复制到目的端（创建镜像）：

```
$ rsync -av source_path destination_path
```

其中 -a 表示要进行归档，-v (verbose) 表示在 stdout 上打印出细节信息或进度。

上面的命令会以递归的方式将所有的文件从源路径复制到目的路径。我们可以指定本地路径，也可以指定远端路径。

路径格式可以是 /home/slynux/data、slynux@192.168.0.6:/home/backups/data 等。

/home/slynux/data 指定了执行 rsync 命令主机上的绝对路径。slynux@192.168.0.6:/home/backups/data 指定了路径 /home/backups/data 位于 IP 地址为 192.168.0.6 的主机上，并以用户 slynux 的身份登录该主机。

要将数据备份到远程服务器或主机，可以使用：

```
$ rsync -av source_dir username@host:PATH
```

如果需要在源端建立一份镜像，只需要定期运行同样的 rsync 命令即可。它只会对更改过的文件进行复制。

用下面的方法将远程主机上的数据恢复到本地主机：

```
$ rsync -av username@host:PATH destination
```

rsync 命令用 SSH 连接远程主机。用 user@host 这种形式设定远程主机的地址，其中 user 代表用户名，host 代表远程主机的 IP 地址或域名。而 PATH 指定需要复制数据的绝对路径。rsync 通常会询问 SSH 连接的密码。这可以通过使用 SSH 密钥来实现自动化（无需用户输入密码）。

请确保远程主机上安装并运行着 OpenSSH。

在通过网络传输时，对数据进行压缩能够明显改善传输速度。我们可以用rsync的选项 `-z` 指定在网络传输时使用数据压缩。例如：

```
$ rsync -avz source destination
```



就路径格式而言，如果我们在source_path末尾使用`/`，那么rsync会将source_path尾端目录中的所有内容复制到目的端。如果没有使用`/`，rsync会将source_path尾端目录本身复制到目的端。

例如，下面的命令会复制test目录中的内容：

```
$ rsync -av /home/test/ /home/backups
```

下面的命令则会将test目录复制到目的端：

```
$ rsync -av /home/test /home/backups
```



如果在destination_path末尾使用`/`，那么rsync会将来自源端的内容复制到目的端目录中。

如果没有使用`/`，rsync会在目的端路径尾部创建一个同名目录，然后将源端内容复制到这个目录中。

例如：

```
$ rsync -av /home/test /home/backups/
```

该命令将源端（/home/test）的内容复制到已有目录backups中。

```
$ rsync -av /home/test /home/backups
```

该命令在目的端创建目录backups，并将源端（/home/test）的内容复制到此目录中。

6

6.10.2 补充内容

可以通过命令行选项来指定rsync命令的其他功能。让我们看看这些选项。

1. 在归档时用rsync排除部分文件

当对远端内容进行归档时，有些文件并不需要进行更新。我们可以告知rsync将某些文件从此次操作中排除。有两个选项可以用来排除文件：

```
--exclude PATTERN
```

可以通过通配符指定需要排除的文件。例如：

```
$ rsync -avz /home/code/some_code /mnt/disk/backup/code --exclude "*.txt"
```

这条命令不对.txt文件进行备份。

或者我们可以通过一个列表文件指定需要排除的文件。

就用`--exclude-from FILEPATH`。

2. 在更新rsync备份时，删除不存在的文件

我们对文件进行归档，然后将归档文件传输到远端备份。当需要更新备份数据时，再创建一个tar文件并传输到远端。默认情况下，rsync并不会在目的端删除那些在源端已不存在的文件。如果要删除这些已不存在文件，使用rsync的--delete选项：

```
$ rsync -avz SOURCE DESTINATION --delete
```

3. 定期调度备份

你可以创建一个cron任务以定期进行备份。

下面是一个简单的例子：

```
$ crontab -e
```

添加上这么一行：

```
0 */10 * * * rsync -avz /home/code user@IP_ADDRESS:/home/backups
```

上面的crontab将rsync调度为每10个小时运行一次。

/10在crontab语法中的小时位(hour position)，/10指定每10小时执行一次备份。如果/10出现在分钟位(minutes position)，那就是每10分钟执行一次备份。

请参阅9.8节了解如何配置crontab。

6.11 用 Git 备份版本控制

人们在备份数据时会采取不同的策略。不过，比起将整个源目录复制到目的端的备份目录中，并使用日期或时间作为版本号，利用差异备份(Differential backup)要更有效。前者只会白白浪费存储空间。而我们只需要复制那些在备份之后发生变化的文件。这被称为增量备份(incremental backup)。我们可以用rsync这类工具手动创建增量备份。不过恢复这种备份可不是件易事。维护和恢复变更的最好方法是使用版本控制系统。由于代码变更频繁，版本控制系统多用于软件开发与代码维护。Git(GNU it)是有名气也是最高效的版本控制系统。我们可以在非编程环境下用Git备份普通文件。通过你使用的发行版的软件包管理器来安装Git。它可是Linus Trovalds编写的哦。

6.11.1 预备知识

下面是有待处理的问题：

我们有一个目录，里面包含了多个文件和子目录。我们需要跟踪目录内容的变更并为其备份。如果数据受损或丢失，我们必须能够恢复数据之前的备份。我们需要将数据定期备份到远程主机上。同样需要在同一台主机(本地主机)的不同位置进行备份。让我们看看如何用Git搞定它。

6.11.2 实战演练

进入需要备份的目录：

```
$ cd /home/data/source
```

使其成为需要跟踪的源目录。

设置并初始化远端备份目录。在远程主机中创建备份目录：

```
$ mkdir -p /home/backups/backup.git
```

```
$ cd /home/backups/backup.git
```

```
$ git init --bare
```

在源主机中执行下列步骤。

(1) 在源主机中添加用户详细信息：

```
$ git config --global user.name "Sarath Lakshman"
```

#将用户名设置为 "Sarath Lakshman"

```
$ git config --global user.email slynux@slynux.com
```

#将email设置为slynux@slynux.com

初始化主机中需要进行备份的源目录。在源目录中执行下列命令：

```
$ git init
```

Initialized empty Git repository in /home/backups/backup.git/

初始化git仓库

```
$ git commit --allow-empty -am "Init"
```

[master (root-commit) b595488] Init

(2) 在源目录中执行下列命令添加远程git目录并同步备份：

```
$ git remote add origin user@remotehost:/home/backups/backup.git
```

```
$ git push origin master
```

Counting objects: 2, done.

Writing objects: 100% (2/2), 153 bytes, done.

Total 2 (delta 0), reused 0 (delta 0)

To user@remotehost:/home/backups/backup.git

* [new branch] master -> master

(3) 为Git跟踪 (Git tracking) 添加或删除文件

下面的命令将当前目录下的所有文件和文件夹添加到备份列表中：

```
$ git add *
```

我们可以有条件地添加某些文件到备份列表中：

```
$ git add *.txt
```

```
$ git add *.py
```

我们可以删除不需要跟踪的文件和文件夹：

```
$ git rm file
```

也可以是文件，甚至是通配符：

```
$ git rm *.txt
```

(4) 检查点或创建备份点

用下列命令利用一个消息来标记备份的检查点 (check point):

```
$ git commit -m "Commit Message"
```

我们需要定期更新远端的备份, 因此需要设置一个cron任务 (例如, 每5个小时进行一次备份)。

用下面的内容创建一个crontab项:

```
0 */5 * * * /home/data/backup.sh
```

创建脚本 /home/data/backup.sh:

```
#!/bin/ bash
cd /home/data/source
git add .
git commit -am "Commit - @ $(date)"
git push
```

现在我们就算是完成了备份系统的设置。

(5) 用Git恢复数据

查看所有的备份版本, 可以使用:

```
$ git log
```

忽略最近的任何变更, 将当前目录更新到最新的备份版本。

□ 要返回任何之前的状态或版本, 要查看一个由32位16进制串组成的提交ID (commit ID)。

通过git checkout来使用提交ID。

□ 如果是提交ID 3131f9661ec1739f72c213ec5769bc0abefa85a9的话, 那就是:

```
$ git checkout 3131f9661ec1739f72c213ec5769bc0abefa85a9
```

```
$ git commit -am "Restore @ $(date) commit ID:"
```

```
3131f9661ec1739f72c213ec5769bc0abefa85a9"
```

```
$ git push
```

□ 要再次查看版本细节信息, 可以使用:

```
$ git log
```

如果工作目录由于某些原因受到了损坏, 我们需要用远端的备份来修复目录。

按照下面的方法从远端备份中重建损坏的内容:

```
$ git clone user@remotehost:/home/backups/backup.git
```

这条命令将创建一个包含全部内容的目录。

6.12 用 dd 克隆磁盘

同硬盘和分区打交道时, 我们可能需要创建所有分区的副本或备份, 而不仅仅是复制内容 (不仅是各个硬盘分区, 而且包括引导记录、分区表等信息)。这时, 我们就可以使用dd命令, 它能够用于克隆任何类型的磁盘, 如硬盘、闪存、CD、DVD以及软盘等。

6.12.1 预备知识

dd命令原意是Data Definiton。由于不正确的使用会导致数据丢失，因此它获得了Data Destroyer这一谑称。使用dd的时候，要留意参数的顺序。错误的参数会损毁全部数据。dd基本上算是一个比特流复制器（bitstream duplicator），它可以将来自磁盘的比特流写入文件，也可以将来自文件的比特流写入磁盘。让我们来看看dd的用法。

6.12.2 实战演练

dd的语法如下：

```
$ dd if=SOURCE of=TARGET bs=BLOCK_SIZE count=COUNT
```

其中：

- if代表输入文件或输入设备路径；
- of代表目标文件或目标设备路径；
- bs代表块大小（通常以2的幂数形式给出，如512、1024、2048等）。COUNT是需要复制的块数（整数）。

需要复制的字节总数 = 块大小 * COUNT

bs和count都是可选的。

通过指定COUNT，我们可以限制从输入文件复制到目标的字节数。如果不指定COUNT，dd会对输入文件进行复制，直到遇见文件结束标志（EOF）为止。

要将一个分区复制到文件中，可以使用：

```
# dd if=/dev/sda1 of=sda1_partition.img
```

该命令中的 /dev/sda1是该分区的设备路径。

用备份恢复分区：

```
# dd if=sda1_partition.img of=/dev/sda1
```

要留心if和of参数。错误使用会造成数据丢失。

将设备路径 /dev/sda1修改成适合的路径，就可以复制或恢复任何磁盘。

如果要永久性删除一个分区中的所有数据，我们可以用dd向该分区中写入0值，命令如下：

```
# dd if=/dev/zero of=/dev/sda1
```

/dev/zero是一个字符设备。它总是返回字符 '\0'。

在容量相同的硬盘间进行克隆：

```
# dd if=/dev/sda of=/dev/sdb
```

其中，/dev/sdb是第二个硬盘。

要制作CD ROM的镜像（ISO文件），方法如下：

```
# dd if=/dev/cdrom of=cdrom.iso
```

6.12.3 补充内容

如果将文件系统创建在一个由dd生成的文件中,那么我们可以将它挂载到挂载点上。下面让我们看看使用方法。

挂载镜像文件

用环回(loopback)的方法可以将任何由dd生成的文件镜像进行挂载,这是用mount命令的-o loop实现的。

```
# mkdir /mnt/mount_point  
# mount -o loop file.img /mnt/mount_point
```

这样,我们就可以通过 /mnt/mount_point来访问镜像文件的内容了。

6.12.4 参考

3.12节讲解了如何利用dd从CD中创建一个ISO文件。



本章内容

- 联网知识入门
- 使用ping
- 列出网络上所有的活动主机
- 通过网络传输文件
- 用脚本设置以太网与无线LAN
- 用SSH实现无密码自动登录
- 用SSH在远程主机上运行命令
- 在本地挂载点上挂载远程驱动器
- 在网络上发送多播式窗口消息
- 网络流量与端口分析

7.1 简介

联网就是通过网络将主机进行互联并采用不同的规范配置网络上的节点。我们以TCP/IP作为网络栈，所有的操作都是基于它进行的。网络是计算机系统中重要的部分。连接在网络上的每个节点都分配了一个用作标识作用的独一的IP地址。有很多联网参数，如子网掩码、路由、端口和DNS等，我们需要对这些知识有一个基本的理解。

一些使用网络的应用通过打开并连接到防火墙端口进行运作。每一个应用都可以提供数据传输、远程shell登录等服务。有些管理任务可以通过网络进行。而shell脚本也能用来配置网络节点、测试主机是否可用、自动化执行远程主机命令等。本章着重介绍网络相关的工具和命令，以及如何用它们解决各种问题。

7.2 联网知识入门

在深入学习与联网相关的攻略之前，有必要简单了解一下网络设置、相关术语以及用于分配IP地址、添加路由等命令。这则攻略会从头开始介绍GNU/Linux中用于联网的命令及用法。

7.2.1 新手上路

网络上每一个节点都需要分配多个参数才能够与其他主机顺利实现互联。这些参数包括子网掩码、网关、路由、端口、DNS等。

这则攻略将介绍命令ifconfig、route、nslookup以及host。

7.2.2 实战演练

网络接口用来连接网络。通常在类UNIX操作系统环境中,网络接口遵循eth0、eth1这种命名惯例。还有一些其他的接口,如usb0、wlan0等,分别对应USB网络接口、无线LAN等网络。

ifconfig命令用于显示网络接口、子网掩码等详细信息。

ifconfig位于/sbin/ifconfig。当输入该命令时,一些GNU/Linux发行版会显示错误“command not found”。这是因为/sbin没有包含在用户的PATH环境变量中。Bash会在PATH变量指定的目录中查找输入的命令。

在Debian中,ifconfig默认是不可用的,因为/sbin并不包括在PATH中。

/sbin/ifconfig是一个绝对路径,所以要用这个绝对路径(也就是/sbin/ifconfig)运行ifconfig。在每个系统中,默认都有一个称之为环回接口的lo,这个接口指向当前主机本身。例如:

```
$ ifconfig
lo        Link encap:Local Loopback
inet addr:127.0.0.1 Mask:255.0.0.0
inet6addr: ::1/128 Scope:Host
          UP LOOPBACK RUNNING MTU:16436 Metric:1
          RX packets:6078 errors:0 dropped:0 overruns:0 frame:0
          TX packets:6078 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:0
          RX bytes:634520 (634.5 KB) TX bytes:634520 (634.5 KB)

wlan0     Link encap:EthernetHWaddr 00:1c:bf:87:25:d2
inet addr:192.168.0.82 Bcast:192.168.3.255 Mask:255.255.252.0
inet6addr: fe80::21c:bfff:fe87:25d2/64 Scope:Link
          UP BROADCAST RUNNING MULTICAST MTU:1500 Metric:1
          RX packets:420917 errors:0 dropped:0 overruns:0 frame:0
          TX packets:86820 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:98027420 (98.0 MB) TX bytes:22602672 (22.6 MB)
```

ifconfig输出的最左边一列是网络接口名,右边的若干列显示对应的网络接口的详细信息。

7.2.3 补充内容

还有一些其他的使用频繁的命令,用于查询及配置网络。让我们看看这些重要的命令以及它们的用法。

1. 打印网络接口列表

这个单行命令序列可以打印系统可用的网络接口列表。

```
$ ifconfig | cut -c-10 | tr -d ' ' | tr -s '\n'
lo
wlan0
```

ifconfig输出到前10个字符是被保留用于打印网络接口的名称,因此我们用cut命令提取每一行的前10个字符。tr -d ' ' 删除每一行的所有空格。用tr -s '\n'压缩重复的换行符生成接口名称列表。

2. IP地址的分配与显示

ifconfig会显示系统中所有可用网络接口的详细信息。我们可以限制它只显示某个特定的接口信息：

```
$ ifconfig iface_name
```

例如：

```
$ ifconfig wlan0
wlan0      Link encap:Ethernet  HWaddr 00:1c:bf:87:25:d2
inet addr:192.168.0.82  Bcast:192.168.3.255
           Mask:255.255.252.0
```

在上面命令输出中，我们感兴趣的是IP地址、广播地址、硬件地址和子网掩码。它们分别为：

□ HWaddr 00:1c:bf:87:25:d2是硬件地址（MAC地址），

□ inet addr:192.168.0.82是IP地址，

□ Bcast:192.168.3.255是广播地址，

□ Mask:255.255.252.0是子网掩码。

在编写某些脚本时，我们可能需要在脚本中提取某些地址来做进一步的处理。

常见的是提取IP地址。要从ifconfig输出中提取IP地址，可以使用：

```
$ ifconfig wlan0 | egrep -o "inet addr:[^ ]*" | grep -o "[0-9.]*"
192.168.0.82
```

第一条命令egrep -o "inet addr:[^]*" 会打印出inet addr:192.168.0.82。

样式以inet addr:作为起始，以非空格字符序列（由[^]*指定）作为结束。在接下来的管道操作中，打印出数字与点号（.）的组合。

要设置网络接口的IP地址，可以使用：

```
# ifconfig wlan0 192.168.0.80
```

你需要以超级用户的身份来运行上面的命令。192.168.0.80是要被设置的IP地址。

设置IP地址的子网掩码：

```
# ifconfig wlan0 192.168.0.80 netmask 255.255.252.0
```

3. 硬件地址（MAC地址）欺骗

在某些情况下，需要利用硬件地址对网络上的计算机进行认证或过滤，对此我们可以使用硬件地址欺骗。硬件地址在ifconfig输出中是以HWaddr 00:1c:bf:87:25:d2的形式出现的。

我们能够按照下面的方法在软件层面上进行硬件地址欺骗：

```
# ifconfig eth0 hw ether 00:1c:bf:87:25:d5
```

在上面的命令中，00:1c:bf:87:25:d5是分配的新MAC地址。

MAC认证服务提供商为每一台主机提供Internet访问，如果我们需要通过他们才能上网，那么上面这招就很管用。

4. 名字服务器与DNS

Internet最根本的寻址方案是IP地址（点分十进制形式，例如202.11.32.75）。然而，Internet

上的资源（比如网站）是通过被称为URL或域名的ASCII字符组合来访问的，例如google.com就是一个域名。它实际上和IP地址是对应关系。在浏览器中输入IP地址同样可以访问www.google.com。

这种利用符号名对IP地址进行抽象的技术就被称为DNS（Domain Name Service，域名服务）。当我们输入google.com，我们的网络所配置的DNS服务器会将域名解析为对应的IP地址。在本地网络中，我们可以设置本地DNS使用主机名命名本地网络上的主机。

分配给当前系统的名字服务器可以通过读取/etc/resolv.conf查看。例如：

```
$ cat /etc/resolv.conf
nameserver 8.8.8.8
```

我们可以像下面这样手动添加名字服务器：

```
# echo nameserver IP_ADDRESS >> /etc/resolv.conf
```

我们该如何获得域名所对应的IP地址呢？

获取IP地址最简单的方法就是ping给定的域名，然后查看回应信息。例如：

```
$ ping google.com
PING google.com (64.233.181.106) 56(84) bytes of data.
#64.233.181.106是google.com对应的IP地址。
```

一个域名可以分配多个IP地址。对于这种情况，DNS服务器只会返回其中一个地址。要想获取分配给域名的所有IP地址，得使用DNS查找工具。

5. DNS查找

有多种的命令行DNS查找工具，这些工具会向DNS服务器请求IP地址解析。而host和nslookup就是两个DNS查找工具。

当执行host时，它会列出某个域名所有的IP地址。nslookup是一个类似于host的命令，它用于查询DNS相关的细节信息以及名字解析。例如：

```
$ host google.com
google.com has address 64.233.181.105
google.com has address 64.233.181.99
google.com has address 64.233.181.147
google.com has address 64.233.181.106
google.com has address 64.233.181.103
google.com has address 64.233.181.104
```

host也可以列出DNS资源记录（DNS resource record），比如MX（Mail Exchange）：

```
$ nslookup google.com
Server:      8.8.8.8
Address:     8.8.8.8#53

Non-authoritative answer:
Name: google.com
Address: 64.233.181.105
Name: google.com
Address: 64.233.181.99
Name: google.com
```

```
Address: 64.233.181.147
Name: google.com
Address: 64.233.181.106
Name: google.com
Address: 64.233.181.103
Name: google.com
Address: 64.233.181.104
```

```
Server: 8.8.8.8
```

上面最后一行对应着用于DNS解析的默认名字服务器。

如果不使用DNS服务器，也可以为IP地址解析添加符号名，这只需要向文件/etc/hosts中添加条目即可。

用下面的语法添加条目：

```
# echo IP_ADDRESS symbolic_name >> /etc/hosts
```

例如：

```
# echo 192.168.0.9 backupserver.com >> /etc/hosts
```

添加了条目之后，任何时候解析backupserver.com，都会返回192.168.0.9。

6. 设置默认网关，显示路由表信息

如果一个本地网络连接到了另一个网络，需要安排一些主机或网络节点，通过它们实现网络间的互联。目的地为外部网络的IP分组应该转发给这些连接外部网络的节点主机。这些能够向外部网络转发分组的特殊节点主机被称为网关。我们要为每个节点设置网关，使之能够连接到外部网络。

操作系统维护着一个被称为路由表（routing table）的表格，它包含了关于分组如何转发以及通过网络中的哪些节点转发的信息。可以用下面的方法显示路由表：

```
$ route
Kernel IP routing table
Destination    Gateway         Genmask         Flags Metric Ref  UseIface
192.168.0.0    *               255.255.252.0   U        2      0    wlan0
link-local     *               255.255.0.0     U       1000   0    wlan0
default        p4.local        0.0.0.0         UG        0      0    wlan0
```

也可以使用：

```
$ route -n
Kernel IP routing table
Destination    Gateway         Genmask         Flags Metric Ref  Use  Iface
192.168.0.0    0.0.0.0         255.255.252.0   U        2      0    0    wlan0
169.254.0.0    0.0.0.0         255.255.0.0     U       1000   0    0    wlan0
0.0.0.0        192.168.0.4     0.0.0.0         UG        0      0    0    wlan0
```

-n指定以数字形式显示地址。如果使用-n，route会以数字形式的IP地址显示每一个条目；否则，在IP地址对应的DNS条目可用的情况下，它会显示符号形式的主机名。

设置默认网关：

```
# route add default gw IP_ADDRESS INTERFACE_NAME
```

例如:

```
# route add default gw 192.168.0.1 wlan0
```

7. Traceroute

当应用程序通过Internet请求服务时,服务器可能位于远端,两者之间可能通过多个网关或设备节点相连。分组要穿过这些网关才能到达目的地。有一个很有意思的命令traceroute,它可以显示分组途径的所有网关的地址。traceroute信息可以帮助我们搞明白分组到达目的地需要经过多少跳(hop)。中途的网关或路由器的数量给出了一个测量网络上两个节点直接距离的度量(metric)。traceroute的输出如下:

```
$ traceroute google.com
traceroute to google.com (74.125.77.104), 30 hops max, 60 byte packets
1 gw-c6509.lxb.as5577.net (195.26.4.1) 0.313 ms 0.371 ms 0.457 ms
2 40g.lxb-fra.as5577.net (83.243.12.2) 4.684 ms 4.754 ms 4.823 ms
3 de-cix10.net.google.com (80.81.192.108) 5.312 ms 5.348 ms 5.327 ms
4 209.85.255.170 (209.85.255.170) 5.816 ms 5.791 ms 209.85.255.172
  (209.85.255.172) 5.678 ms
5 209.85.250.140 (209.85.250.140) 10.126 ms 9.867 ms 10.754 ms
6 64.233.175.246 (64.233.175.246) 12.940 ms 72.14.233.114
  (72.14.233.114) 13.736 ms 13.803 ms
7 72.14.239.199 (72.14.239.199) 14.618 ms 209.85.255.166
  (209.85.255.166) 12.755 ms 209.85.255.143 (209.85.255.143) 13.803 ms
8 209.85.255.98 (209.85.255.98) 22.625 ms 209.85.255.110
  (209.85.255.110) 14.122 ms
*
9 ew-in-f104.1e100.net (74.125.77.104) 13.061 ms 13.256 ms 13.484 ms
```

7.2.4 参考

- 1.3节讲解了PATH变量。
- 4.3节讲解了grep命令。

7.3 使用 ping

ping是每位用户都应该首先了解的最基础的网络命令。绝大多数操作系统上都包含该命令。ping也是一个验证网络上两台主机连通性的诊断工具,能够用来找出网络上的活动主机。让我们看看这个命令的用法。

7.3.1 实战演练

为了检查网络上两台主机之间的连通性,ping命令使用互联网控制消息协议(Internet Control Message Protocol, ICMP)的echo分组。当这些echo分组发送到某个主机时,如果分组能够送达且该主机为活动主机,那么它就会发送一条回应。

检查某个主机是否可以到达:

```
$ ping ADDRESS
```

ADDRESS可以是主机名、域名或者IP地址。

ping会连续发送分组，回应信息将被打印在终端上。用Ctrl+C来停止ping命令。

例如：

□ 如果主机可以到达，那么输出信息如下所示。

```
$ ping 192.168.0.1
PING 192.168.0.1 (192.168.0.1) 56(84) bytes of data.
64 bytes from 192.168.0.1: icmp_seq=1 ttl=64 time=1.44 ms
^C
--- 192.168.0.1 ping statistics ---
1 packets transmitted, 1 received, 0% packet loss, time 0ms
rtt min/avg/max/mdev = 1.440/1.440/1.440/0.000 ms

$ ping google.com
PING google.com (209.85.153.104) 56(84) bytes of data.
64 bytes from bom01s01-in-f104.1e100.net (209.85.153.104): icmp_
seq=1 ttl=53 time=123 ms
^C
--- google.com ping statistics ---
1 packets transmitted, 1 received, 0% packet loss, time 0ms
rtt min/avg/max/mdev = 123.388/123.388/123.388/0.000 ms
```

□ 如果主机不可到达，则输出如下所示。

```
$ ping 192.168.0.99
PING 192.168.0.99 (192.168.0.99) 56(84) bytes of data.
From 192.168.0.82 icmp_seq=1 Destination Host Unreachable
From 192.168.0.82 icmp_seq=2 Destination Host Unreachable
```

一旦主机不可到达，ping返回错误信息“Destination Host Unreachable”。

7.3.2 补充内容

除了检查网络上两点之间的连通性，ping命令还可以通过其他选项来获取有用信息。让我们看看ping的其他选项。

1. 往返时间

ping命令可以用来得出网络上两台主机之间的往返时间（Round Trip Time, RTT）。RTT是分组从源主机到目的主机的往返时间。RTT的单位是毫秒，可以从ping命令中获知。例如：

```
--- google.com ping statistics ---
5 packets transmitted, 5 received, 0% packet loss, time 4000ms
rtt min/avg/max/mdev = 118.012/206.630/347.186/77.713 ms
```

其中，最小的RTT是118.012ms，平均RTT是206.630ms，最大RTT是347.186ms，ping输出中的mdev（77.713ms）代表平均偏差。

2. 限制发送的分组数量

ping命令会不停地发送echo分组，并等待回复，直到按下Ctrl+C。不过，我们可以用选项 -c

限制所发送的echo分组的数量。

用法如下：

```
-c COUNT
```

例如：

```
$ ping 192.168.0.1 -c 2
PING 192.168.0.1 (192.168.0.1) 56(84) bytes of data.
64 bytes from 192.168.0.1: icmp_seq=1 ttl=64 time=4.02 ms
64 bytes from 192.168.0.1: icmp_seq=2 ttl=64 time=1.03 ms
--- 192.168.0.1 ping statistics ---
2 packets transmitted, 2 received, 0% packet loss, time 1001ms
rtt min/avg/max/mdev = 1.039/2.533/4.028/1.495 ms
```

在上面的例子中，ping命令发送了2个echo分组后就停止发送。

如果我们需要通过脚本ping一组IP地址来检查主机的状态，那么这个技巧就能派上用场。

3. ping命令的返回状态

ping命令如果执行顺利，会返回退出状态0；否则，返回非0。执行顺利意味着目的主机能够到达，否则意味着目的主机不可到达。

返回状态可以通过下面的方法轻松获得：

```
$ ping ADDRESS -c2
if [ $? -eq 0 ];
then
    echo Successful ;
else
    echo Failure
fi
```

7.4 列出网络上所有的活动主机

当我们涉及大型局域网时，可能需要检查网络上的其他主机是否处于活动状态。一台非活动主机可能有两种情况：要么是没有开机，要么是网络连接有问题。借助shell脚本，我们可以轻易找出并报告网络上的哪一台主机处于活动状态。让我们看看这是如何实现的。

7.4.1 新手上路

在这则攻略中，我们采用两种方法。第一种方法使用ping，第二种方法使用fping。fping并没有默认包含在Linux发行版中。你可以用软件包管理器手动安装。

7.4.2 实战演练

让我们看看可以找出网络上所有活动主机的脚本，以及实现这一目标的其他方法。

□ 方法1

我们可以用ping命令编写脚本来查询一组IP地址同时检查它们是否处于活动状态：

```
#!/bin/bash
#文件名: ping.sh
#用途: 根据你的网络配置对网络地址192.168.0进行修改.
```

```
for ip in 192.168.0.{1..255} ;
do
    ping $ip -c 2 &> /dev/null ;

    if [ $? -eq 0 ] ;
    then
        echo $ip is alive
    fi
done
```

输出如下:

```
$ ./ping.sh
192.168.0.1 is alive
192.168.0.90 is alive
```

□ 方法2

我们可以用已有的命令行工具来查询网络上的主机状态:

```
$ fping -a 192.160.1/24 -g 2> /dev/null
192.168.0.1
192.168.0.90
```

或者, 使用:

```
$ fping -a 192.168.0.1 192.168.0.255 -g
```

7.4.3 工作原理

在方法1中, 我们用ping命令找出网络上的活动主机。我们用了一个for循环对一组IP地址进行迭代。这组IP地址依照192.168.0.{1..255}格式生成。像{start..end}这种记法会由shell对其进行扩展并生成一组IP地址, 例如, 192.168.0.1、192.168.0.2、192.168.0.3等依次类推, 直到192.168.0.255。

ping \$ip -c 2 &> /dev/null会在每次循环中ping对应的IP地址。-c 2将发送的echo分组数量限制为2。&> /dev/null用于将stderr和stdout重定向到/dev/null, 使它们不会在终端上打印出来。我们用\$?获取退出状态。如果顺利退出, 退出状态为0, 否则, 为非0。因此能够ping通的IP地址就被打印出来。我们也可以打印出不能ping通的IP地址, 从而给出一个无法抵达的IP地址列表。



这里为你准备了一个练习。不是在脚本中使用硬编码形式的IP地址范围, 而是修改脚本使其能够从文件或stdin中读取一组IP地址。

在这个脚本中，每一次ping都是依次执行的。即使所有的IP地址都是彼此独立，由于编写的是顺序式程序（sequential program），ping命令也只能按顺序执行。每执行一次ping命令，都要经历一段延迟：发送两个echo分组，并接收回应或等待回应超时。

要是处理255个地址的话，这段延迟的时间可就不短了。我们可以使用并行方式来加速所有ping命令的执行。这个脚本的核心部分是循环体。要使ping命令可以并行执行，可将循环体放入（）&。将命令块放入（），使其中的命令可作为子shell来执行，而&可以使之脱离当前线程，在后台继续运行。例如：

```
(
    ping $ip -c2 &> /dev/null ;
    if [ $? -eq 0 ];
    then
        echo $ip is alive
    fi
)&

wait
```

for循环体执行多个后台进程，然后结束循环并终止脚本。要想等所有子进程结束之后再终止脚本，我们得使用wait命令。将wait放在脚本最后，它就会一直等到所有的子进程全部结束。



wait命令使脚本只有在所有的子进程或后台进程全部终止或完成之后才能结束。

请阅读本书所提供的脚本fast_ping.sh。

方法2使用了另一个命令fping。它同时可以ping一组IP地址，而且响应速度非常快。fping的选项如下：

- 选项 -a 指定打印出所有活动主机的IP地址。
- 选项 -u 指定打印出所有无法到达的主机。
- 选项 -g 指定从写作IP/mask的“斜线-子网掩码”记法或者起止IP地址记法中生成IP地址范围。

```
$ fping -a 192.160.1/24 -g
```

或者

```
$ fping -a 192.160.1 192.168.0.255 -g
```

- 2>/dev/null 将由于主机无法到达所产生的错误信息打印到null设备。

也可以采用命令行参数的方式手动指定一组IP地址，或者从stdin中接收。例如：

```
$ fping -a 192.168.0.1 192.168.0.5 192.168.0.6
```

```
# 将IP地址作为参数传递
```

```
$ fping -a <ip.list
```

```
# 从文件中传递一组IP地址
```

7.4.4 补充内容

`fping`命令可以用来查询网络中的DNS数据。来看看它的具体用法。

用**`fping`**进行DNS查找

`fping`有一个选项`-d`，它通过对每一个echo回应进行DNS查询来返回主机名。使用该选项可以在ping的回应信息中打印出主机名而非IP地址。

```
$ cat ip.list
192.168.0.86
192.168.0.9
192.168.0.6

$ fping -a -d 2>/dev/null <ip.list
www.local
dns.local
```

7.4.5 参考

- 1.5节讲解了数据重定向。
- 1.15节讲解了数字比较。

7.5 传输文件

计算机联网的主要目的就是资源共享。在资源共享方面，使用最多的是文件共享。有很多不同的方法可以在网络节点上传输文件。这则攻略就讨论了如何用常见的协议FTP、SFTP、RSYNC和SCP传输文件。

7.5.1 新手上路

用来在网络上传输文件的命令多数都已默认包含在安装好的Linux中。通过FTP传输文件可以使用`lftp`命令，通过SSH传输文件可以使用`sftp`，RSYNC使用SSH与`rsync`命令，并借助`scp`通过SSH进行传输。

7.5.2 实战演练

文件传输协议（File Transfer Protocol, FTP）是一个古老的用于网络主机之间传输文件的文件传输协议。我们可以用`lftp`命令访问FTP服务器以便传输文件。它使用端口21。只有远程主机上安装有FTP服务器才能使用FTP。很多公共网站都是用FTP共享文件。

要连接FTP服务器传输文件，可以使用：

```
$ lftp username@ftphost
```

它会提示你输入密码，然后显示一个像下面那样的登录提示符：

```
lftp username@ftphost:~>
```

你可以在提示符后输入命令，如下所示。

- 用cd directory改变目录。
- 用lcd改变本地主机的目录。
- 用mkdir创建目录。
- 用get filename下载文件：

```
lftp username@ftphost:~> get filename
```

- 用put filename从当前目录上传文件：

```
lftp username@ftphost:~> put filename
```

- 用quit退出lftp会话。

lftp提示符支持命令自动补全。

7.5.3 补充内容

让我们看看其他可用于网络文件传输的技术及命令。

1. FTP自动传输

ftp是一个基于FTP文件传输的命令。相比较而言，lftp的用法更灵活。lftp和ftp为用户启动一个交互式会话（通过显示消息提示用户输入）。如果我们不使用交互模式，而是希望进行自动文件传输，又该怎么做呢？下面的脚本可以用来实现FTP自动传输。

```
#!/bin/bash
#文件名: ftp.sh
#用途: 自动化FTP传输
HOST='domain.com'
USER='foo'
PASSWD='password'
ftp -i -n $HOST <<EOF
user ${USER} ${PASSWD}
binary
cd /home/slynux
puttestfile.jpg
getserverfile.jpg
quit
EOF
```

上面的脚本包含下列结构：

```
<<EOF
DATA
EOF
```

这是用来通过stdin向FTP命令发送数据。1.5节中已讲解了重定向到stdin的各种方法。

在示例脚本中，ftp的选项-i关闭用户的交互会话，user \${USER} \${PASSWD}设置用户名和密码，binary将文件模式设置为二进制模式。

2. SFTP

SFTP (Secure FTP, 安全FTP) 是一个类似于FTP的文件传输系统, 它运行在SSH连接之上。SFTP利用SSH连接模拟FTP接口。它不需要远端运行FTP服务器来执行文件传输, 但必须安装并运行OpenSSH服务器。SFTP是一个交互式命令, 提供了命令提示符。

下面的命令用来执行文件传输。对于特定主机、用户和密码的自动化FTP会话来说, 余下的命令都是一样的。

```
cd /home/slynux
put testfile.jpg
get serverfile.jpg
```

运行sftp:

```
$ sftp user@domainname
```

和lftp类似, 输入quit命令可以退出sftp会话。

SSH服务器有时候并不在默认的端口22上运行。如果它在其他端口运行, 我们可以在sftp中用 `-oPort=PORTNO` 来指定端口号。

例如:

```
$ sftp -oPort=422 user@slynux.org
```



`-oPort`应该作为sftp命令的第一个参数。

3. RSYNC

作为一款重要的命令行工具, `rsync`广泛用于网络文件与系统快照的备份。6.10节详细讲解了`rsync`的用法。

4. SCP

SCP (Secure Copy, 安全复制) 是一项比传统远程复制工具`rcp`更安全的文件复制技术。文件都是通过SSH加密通道进行传输的。我们可以像下面这样轻松地将文件传输到远程主机:

```
$ scp filename user@remotehost:/home/path
```

该命令会提示你输入密码。我们可以用SSH自动登录功能来免于输入密码。7.7节讲解了SSH自动登录。

因此, 用`scp`传输文件无需特定的脚本。一旦实现了SSH自动登录, `scp`就可以直接执行, 而不需要再提示输入密码。

命令中的`remotehost`可以使IP地址或域名。 `scp`命令的格式是:

```
$ scp SOURCE DESTINATION
```

SOURCE或DESTINATION可以采用形如`username@localhost:/path`的格式。例如:

```
$ scp user@remotehost:/home/path/filename filename
```

上面的命令将远程主机中的文件复制到当前目录, 并采用给定的文件名。

如果SSH没有运行在端口22, 使用 `-oPort`, 并采用和 `sftp` 相同的语法。

5. 用 `scp` 进行递归复制

使用 `scp` 的 `-r` 选项, 我们可以在两台网络主机之间对文件夹进行递归复制:

```
$ scp -r /home/slynux user@remotehost:/home/backups
# 将目录/home/slynux递归复制到远程主机中
```

`scp` 的 `-p` 选项能够在复制文件的同时保留文件的权限和模式。

7.5.4 参考

1.5节讲解了如何用EOF实现标准输入。

7.6 用脚本设置以太网与无线 LAN

配置以太网很简单, 因为它使用物理线缆, 没有认证之类的特殊要求。但是无线LAN就需要认证了, 比如WEP密钥, 还有所要连接的无线网络的ESSID。让我们看看如何用脚本来连接无线网络和有线网络。

7.6.1 新手上路

我们需要用 `ifconfig` 分配IP地址和子网掩码才能连接上有线网络。对于无线网络来说, 还需要其他工具 (如 `iwconfig` 和 `iwlist`) 来配置更多的参数。

7.6.2 实战演练

要通过有线网络接口连接网络, 执行下面的脚本:

```
#!/bin/bash
#文件名: etherconnect.sh
#用途: 连接以太网

#根据你的设置修改下面的参数
##### PARAMETERS #####

IFACE=eth0
IP_ADDR=192.168.0.5
SUBNET_MASK=255.255.255.0
GW=192.168.0.1
HW_ADDR='00:1c:bf:87:25:d2'
# HW_ADDR 是可选的
#####

if [ $UID -ne 0 ];
then
    echo "Run as root"
    exit 1
```



```

fi

#进行新的配置之前关闭接口
/sbin/ifconfig $IFACE down

if [[ -n $HW_ADDR ]];
then

    /sbin/ifconfig hw ether $HW_ADDR
    echo Spoofed MAC ADDRESS to $HW_ADDR

fi

/sbin/ifconfig $IFACE $IP_ADDR netmask $SUBNET_MASK

route add default gw $GW $IFACE

echo Successfully configured $IFACE

用WEP连接无线LAN的脚本如下:

#!/bin/bash
#文件名: wlan_connect.sh
#用途: 连接无线 LAN

#根据你的设置修改下面的参数
##### PARAMETERS #####
IFACE=wlan0
IP_ADDR=192.168.1.5
SUBNET_MASK=255.255.255.0
GW=192.168.1.1
HW_ADDR='00:1c:bf:87:25:d2'
#如果不想欺骗物理地址, 在行上注释

ESSID="homenet"
WEP_KEY=8b140b20e7
FREQ=2.462G
#####

KEY_PART=""

if [[ -n $WEP_KEY ]];
then
    KEY_PART="key $WEP_KEY"
fi

#进行新的配置之前关闭接口
/sbin/ifconfig $IFACE down

if [ $UID -ne 0 ];
then
    echo "Run as root"
    exit 1;
fi

```



```

if [[ -n $HW_ADDR ]];
then
    /sbin/ifconfig $IFACE hw ether $HW_ADDR
    echo Spoofed MAC ADDRESS to $HW_ADDR
fi

/sbin/iwconfig $IFACE essid $ESSID $KEY_PART freq $FREQ

/sbin/ifconfig $IFACE $IP_ADDR netmask $SUBNET_MASK

route add default gw $GW $IFACE

echo Successfully configured $IFACE

```

7.6.3 工作原理

命令ifconfig、iwconfig和route必须以超级用户身份来运行，因此在脚本一开始要检查是否为超级用户。

连接以太网的这个脚本相当简单明了，例子中涉及的概念在7.2节中都已讲解过。让我们看看用来连接无线LAN的命令。

无线LAN需要essid、密钥、网络频率等参数。essid是我们想要连接的无线网络的名称。一些使用WEP (Wired Equivalent Protocol, 有线等效保密协议) 的网络需要用WEP密钥进行验证，有些网络则不需要这一步。WEP密钥通常是一个10位十六进制数口令。频率则是分配给特定网络的。iwconfig命令用来为无线网卡配置适合的无线网络、WEP密钥以及频率。

我们可以用iwlist工具扫描并列出可用的无线网络。用下面的命令进行扫描：

```

# iwlist scan
wlan0      Scan completed :
           Cell 01 - Address: 00:12:17:7B:1C:65
                Channel:11
                Frequency:2.462 GHz (Channel 11)
                Quality=33/70  Signal level=-77 dBm
                Encryption key:on
                ESSID:"model-2"

```

可以从扫描结果中的Frequency:2.462 GHz (Channel 11)行提取频率参数。

7.6.4 参考

1.15节讲解了字符串比较。

7.7 用SSH实现无密码自动登录

SSH广泛用于脚本自动化。借助SSH，我们可以在远程主机上执行命令并读取输出。SSH使用用户名和密码进行认证。在SSH命令的执行过程中提示输入密码。但是在自动化脚本中，SSH

命令可能在一个循环中执行上百次，每次都提供密码的话，显然不实际。因此，我们需要将登录过程自动化。SSH就包含了一个内建的特性，可以用SSH密钥实现自动登录。这则攻略描述了如何创建SSH密钥并协助实现自动登录。

实战演练

SSH采用基于公钥和基于私钥的加密技术进行自动化认证。认证密钥包含两部分：一个公钥和一个私钥。我们可以通过ssh-keygen命令创建认证密钥。要想实现自动化认证，公钥必须放置在服务器中（将其加入文件 ~/.ssh/authorized_keys），与公钥对应的私钥应该放入你用来登录的客户机的 ~/.ssh 目录中。另一些与SSH相关的配置信息（例如，authorized_keys文件的路径与名称）可以通过修改文件 /etc/ssh/sshd_config 进行配置。

设置SSH自动化认证需要两步。

(1) 创建SSH密钥，这需要登录到远程主机。

(2) 将生成的公钥传输到远程主机，并将其加入文件 ~/.ssh/authorized_keys 中。

要创建SSH密钥，输入命令 ssh-keygen，并规定加密算法类型为RSA：

```
$ ssh-keygen -t rsa
Generating public/private rsa key pair.
Enter file in which to save the key (/home/slynux/.ssh/id_rsa):
Created directory '/home/slynux/.ssh'.
Enter passphrase (empty for no passphrase):
Enter same passphrase again:
Your identification has been saved in /home/slynux/.ssh/id_rsa.
Your public key has been saved in /home/slynux/.ssh/id_rsa.pub.
The key fingerprint is:
f7:17:c6:4d:c9:ee:17:00:af:0f:b3:27:a6:9c:0a:05slynux@slynux-laptop
The key's randomart image is:
+--[ RSA 2048 ]-----+
|           .           |
|            o . .      |
|       E      o o .    |
|      ...oo |         |
|       .S .+ .+o.      |
|      . . .+.....    |
|     .+o...           |
|      . + o . .       |
|      ..+             |
+-----+

```

你需要输入一个口令来生成一对公钥和私钥。如果不输入口令的话，也可以生成密钥，但是这种做法是一种不安全的行为。我们可以编写监控脚本，利用自动登录来登入多台主机。对于这种情况，在运行ssh-keygen命令时，不要填入口令，这样就能够避免在脚本运行时向你索要口令了。

现在 ~/.ssh/id_rsa.pub 和 ~/.ssh/id_rsa 已经生成了。id_rsa.pub 是生成的公钥，id_rsa 是生成的私钥。公钥必须添加到远程服务器 ~/.ssh/authorized_keys 文件中，这台服务器也就是我们想从当前主机自动登入的那台服务器。

要添加一个密钥，可以使用：

```
$ ssh USER@REMOTE_HOST "cat >> ~/.ssh/authorized_keys" < ~/.ssh/id_rsa.
pub
Password:
```

上面的命令要设定登录密码。

这样一来，自动登录就设置好了。从现在开始，SSH在运行过程中就不会再提示输入密码。你可以用下面的命令测试：

```
$ ssh USER@REMOTE_HOST uname
Linux
```

它不会提示你输入密码了。

7.8 用SSH在远程主机上运行命令

SSH是一个很有意思的系统管理工具，它能够通过shell登录并控制远程主机。SSH是Secure Shell的缩写。我们可以登录到远程主机上来执行shell命令，就好像是在本地主机上执行这些命令一样。SSH使用加密通道来传输数据。这则攻略介绍了在远程主机上运行命令的各种方法。

7.8.1 新手上路

所有的GNU/Linux发行版都不默认包括SSH，要用软件包管理器安装openssh-server和openssh-client软件包。SSH服务默认在端口22运行。

7.8.2 实战演练

用下面的方法连接运行SSH服务器的远程主机：

```
$ ssh username@remote_host
```

其中

- username是远程主机上的用户，
- remote_host可以是域名或IP地址。

例如：

```
$ ssh mec@192.168.0.1
The authenticity of host '192.168.0.1 (192.168.0.1)' can't be established.
RSA key fingerprint is 2b:b4:90:79:49:0a:f1:b3:8a:db:9f:73:2d:75:d6:f9.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added '192.168.0.1' (RSA) to the list of known hosts.
Password:
```

```
Last login: Fri Sep 3 05:15:21 2010 from 192.168.0.82
mec@proxy-1:~$
```

ssh采用交互方式询问用户密码，一旦认证成功，将会为用户返回一个shell。

SSH服务器默认在端口22运行。不过，有些运行SSH服务的服务器并不在这个端口运行。针对这种情况，用ssh命令的-p port_no指定端口。

要连接运行在端口422之上的SSH服务器，可以使用：

```
$ ssh user@localhost -p 422
```

你可以在远程主机的shell中执行命令。shell是一个交互式工具，用户可以在其中输入命令并执行它。然而在shell脚本环境中，我们不需要交互式的shell，而是需要将任务自动化。我们需要在远程shell中执行若干命令，然后在本地主机上显示或存储输出内容。每次都要输入密码对于自动化脚本来说显然不实际，因此要对SSH进行自动登录配置。

7.7节讲解了SSH命令。

在运行使用了SSH的自动化脚本之前，要确保配置了自动登录。

要想在远程主机中运行命令，并将命令输出显示在本地shell，使用下面的语法：

```
$ ssh user@host 'COMMANDS'
```

例如：

```
$ ssh mec@192.168.0.1 'whoami'
Password:
mec
```

可以输入多条命令，在命令之间以分号进行分隔：

```
$ ssh user@host 'command1 ; command2 ; command3'
```

命令可以通过stdin发送，而命令输出可以通过stdout来获取。

语法如下：

```
$ ssh user@remote_host "COMMANDS" > stdout.txt 2> errors.txt
```

COMMANDS部分应该用引号括起来，以避免分号被本地主机的shell当做定界符。我们也可以通过stdin将包含管道的命令序列传递给SSH命令：

```
$ echo "COMMANDS" | sshuser@remote_host> stdout.txt 2> errors.txt
```

例如：

```
$ ssh mec@192.168.0.1 "echo user: $(whoami);echo OS: $(uname)"
Password:
user: slynux
OS: Linux
```

在这个例子中，远程主机上执行的命令是：

```
echo user: $(whoami);
echo OS: $(uname)
```

写成一般化的形式就是：

```
COMMANDS="command1; command2; command3"
$ ssh user@hostname "$COMMANDS"
```

我们也可以在命令序列中用()子shell操作符传递一个更复杂的子shell。

现在,编写一个基于SSH的shell脚本,用来收集一组远程主机的运行时间(uptime)。运行时间是系统加电运行的时间。uptime命令用来显示系统加电后运行了多久。

假设IP_LIST中的所有系统都有一个用户test。

```
#!/bin/bash
#文件名: uptime.sh
#用途: 系统运行时间监视器

IP_LIST="192.168.0.1 192.168.0.5 192.168.0.9"
USER="test"

for IP in $IP_LIST;
do
    uptime=$(ssh $USER@$IP uptime | awk '{ print $3 }')
    echo $IP uptime: $uptime
done
```

输出应该是:

```
$ ./uptime.sh
192.168.0.1 uptime: 1:50,
192.168.0.5 uptime: 2:15,
192.168.0.9 uptime: 10:15,
```

7.8.3 补充内容

让我们看看ssh命令的另一些选项。

1. SSH的压缩功能

SSH协议也支持对数据进行压缩传输,当带宽有限时,这一功能很方便。用ssh命令的选项-C启用压缩功能:

```
$ ssh -C user@hostname COMMANDS
```

2. 将数据重定向至远程shell命令的stdin

有时候我们需要将一些数据重定向到远程shell命令的stdin。通过下面的例子看看这如何实现:

```
$ echo "text" | ssh user@remote_host 'cat >> list'
```

或者:

```
# 将文件中的数据重定向
$ ssh user@remote_host 'cat >> list' < file
```

cat >> list将通过stdin接收到的数据添加到文件list中。这条命令是在远程主机上执行的,但数据却是从本地主机传递到远程shell的stdin的。

7.8.4 参考

7.7节讲解了如何配置自动登录来执行命令而无需提示输入密码。

7.9 在本地挂载点上挂载远程驱动器

在执行读写数据传输操作时，如果可以通过本地挂载点访问远程主机文件系统，那就再好不过。SSH是网络中最常用的文件传输协议，因此我们利用它和sshfs来实现这一需求。sshfs允许你将远程文件系统挂载到本地挂载点上。让我们看看实现方法。

7.9.1 新手上路

GNU/Linux并不默认包含sshfs。用软件包管理器安装这个工具。sshfs是FUSE文件系统的一个扩展，FUSE允许其支持的操作系统像使用本地文件系统一样挂载各类数据。

7.9.2 实战演练

要将位于远程主机上的文件系统挂载到本地挂载点上，可以使用：

```
# sshfs user@remotehost:/home/path /mnt/mountpoint
Password:
```

在收到提示时输入用户密码。

现在位于远程主机/home/path中的数据就可以通过本地挂载点/mnt/mountpoint来访问了。

完成任务后，可用下面的方法卸载：

```
# umount /mnt/mountpoint
```

7.9.3 参考

7.8节讲解了ssh命令。

7.10 在网络上发送多播式窗口消息

网络管理员通常需要发送消息到网络上的节点。在用户桌面上显示弹出窗有助于提醒用户一些信息。用GUI工具和shell脚本就可以完成这一任务。这则攻略讨论如何向远程主机发送带有定制消息的弹出窗口。

7.10.1 新手上路

可以用zenity来实现一个GUI弹出窗口。zenity是一个脚本化的GUI工具，用来创建包括文本框、输入框等在内的窗口。SSH可以用来连接远程主机上的shell。GUN/Linux发行版不默认包括zenity，可以用软件包管理器自行安装。

7.10.2 实战演练

zenity是一种脚本化对话框创建工具。还有其他的对话框，如gdialog、kdialog、xdialog

等。作为面向GNOME桌面环境的一款工具，zenity的用法非常灵活。

要用zenity创建一个信息框，可以使用：

```
$ zenity --info --text "This is a message"
# 这条命令会显示一个窗口，上面显示文本 "This is a message"
```

zenity能够创建带有输入框、多选输入、单选输入、按钮等元素的窗口。不过这些内容超出了本攻略的讨论范围，请查看zenity的手册页以获取更多的信息。

现在，我们可以使用SSH在远程主机上运行zenity语句。要通过SSH在远程主机上执行上面那条语句，使用：

```
$ ssh user@remotehost 'zenity --info --text "This is a message"'
```

不过，这会返回一个类似于下面的错误：

```
(zenity:3641): Gtk-WARNING **: cannot open display:
```

这是因为zenity依赖于Xserver。Xserver是一个守护进程，它负责在屏幕上绘制GUI图形元素。精简版的GNU/Linux系统（bare GNU/Linux system）只包含一个文本终端或一些shell提示符。

Xserver用特殊的环境变量DISPLAY跟踪运行在系统中的Xserver实例。

我们可以手动设置DISPLAY=:0来指示Xserver有关其运行实例的情况。

之前的SSH命令可以改写成：

```
$ ssh username@remotehost 'export DISPLAY=:0 ; zenity --info --text "This is a message"'
```

用户一旦登入任何窗口管理器，这条语句就会在remotehost上显示一个弹出窗口。

要以多播的方式向多台远程主机发送弹出窗口，使用下面的脚本：

```
#!/bin/bash
# 文件名: multi_cast_window.sh
# 用途: 以多播形式弹出窗口

IP_LIST="192.168.0.5 192.168.0.3 192.168.0.23"
USER="username"

COMMAND='export DISPLAY=:0 ; zenity --info --text "This is a message" '
for host in $IP_LIST;
do
    ssh $USER@$host "$COMMAND" &
done
```

7.10.3 工作原理

在上面的脚本中，我们有一组需要在其上弹出窗口的IP地址。使用循环对IP地址进行迭代并执行SSH命令。

在SSH语句的末尾有一个 &，这个符号的作用是将SSH语句放置到后台运行。这样做有利于多条SSH语句并行执行。如果不使用 &，ssh命令会启动一个SSH会话，执行zenity，然后等待

用户关闭弹出窗口。除非远程主机上的用户关闭窗口，否则循环中的下一条SSH语句永远都不会被执行。要想摆脱这种必须等待SSH会话结束才能继续执行命令的情况，我们就得使用 `&`。

7.10.4 参考

7.8节讲解了ssh命令。

7.11 网络流量与端口分析

网络端口是网络应用程序必不可少的参数。应用程序在主机上打开端口，然后通过远程主机中打开的端口实现远程通信。出于安全方面的考虑，必须留意系统中打开及关闭的端口。恶意软件和rootkits可能会利用特定的端口及服务运行在系统中，从而使攻击者可以对数据及资源进行未经授权访问。通过获取开放端口列表以及运行在端口之上的服务，我们便可以分析并抵御rootkits的操纵，而且这些信息还有助于我们对其进行清除。开放端口列表不仅能够协助检测恶意软件，而且还能够收集开放端口的相关信息对网络应用程序进行调试。它可以用来分析某些端口连接及端口侦听功能是否正常。这则攻略讨论了各种用于端口分析的工具。

7.11.1 新手上路

很多命令可用来侦听端口以及运行在端口上的服务（例如lsof和netstat）。这些命令都默认包含在GNU/Linux发行版中。

7.11.2 实战演练

要列出系统中的开放端口以及运行在端口上的服务的详细信息，可以使用：

```
$ lsof -i
COMMAND  PID  USER  FD  TYPE  DEVICE  SIZE/OFF  NODE NAME
firefox-b 2261 slynux 78u IPv4  63729      0t0  TCP localhost:47797->localhost:42486 (ESTABLISHED)
firefox-b 2261 slynux 80u IPv4  68270      0t0  TCP slynux-laptop.local:41204->192.168.0.2:3128 (CLOSE_WAIT)
firefox-b 2261 slynux 82u IPv4  68195      0t0  TCP slynux-laptop.local:41197->192.168.0.2:3128 (ESTABLISHED)
ssh       3570 slynux 3u  IPv6  30025      0t0  TCP localhost:39263->localhost:ssh (ESTABLISHED)
ssh       3836 slynux 3u  IPv4  43431      0t0  TCP slynux-laptop.local:40414->honeym.mtveurope.org:422 (ESTABLISHED)
GoogleTel 4022 slynux 12u IPv4  55370      0t0  TCP localhost:42486 (LISTEN)
GoogleTel 4022 slynux 13u IPv4  55379      0t0  TCP localhost:42486->localhost:32955 (ESTABLISHED)
```

lsof的每一项输出都对应对应着一个打开了特定端口的服务。输出的最后一类类似于：

```
slynux-laptop.local:34395->192.168.0.2:3128 (ESTABLISHED)
```

输出中的slynux-laptop.local:34395对应本地主机, 192.168.0.2:3128对应远程主机。

34395是本地主机当前的开放端口, 3128是连接远程主机服务所需要使用的端口。

要列出本地主机当前的开放端口, 可以使用:

```
$ lsdf -l | grep ":[0-9]\|-->" -o | grep "[0-9]\|-->" -o | sort | uniq
```

其中, 第一个grep中使用的正则表达式:[0-9]\|-->用来从lsdf输出中提取主机端口部分(:34395-->); 第二个grep用来提取端口号(数字)。同一个端口可能会有多个连接, 因此相同的端口也许会出现多次。为了保证每个端口只显示一次, 将端口号排序并打印出不重复的部分。

7.11.3 补充内容

让我们看看另一些用来查看开放端口以及网络流量相关信息的工具。

用netstat查看开放端口与服务

netstat是另一个用于网络服务分析的命令。讲解netstat的所有特性将超出这则攻略的范围。我们只看如何用它列出服务与端口号。

用netstat -tnp列出开放端口与服务:

```
$ netstat -tnp
(Not all processes could be identified, non-owned process info
will not be shown, you would have to be root to see it all.)
Active Internet connections (w/o servers)
Proto Recv-Q Send-Q Local Address           Foreign Address         State
PID/Program name
tcp        0      0 192.168.0.82:38163      192.168.0.2:3128       ESTABLISHED
2261/firefox-bin
tcp        0      0 192.168.0.82:38164      192.168.0.2:3128       TIME_
WAIT      -
tcp        0      0 192.168.0.82:40414      193.107.206.24:422     ESTABLISHED
3836/sah
tcp        0      0 127.0.0.1:42486         127.0.0.1:32955        ESTABLISHED
4022/GoogleTalkPlug
tcp        0      0 192.168.0.82:38152      192.168.0.2:3128       ESTABLISHED
2261/firefox-bin
tcp6       0      0 :::1:22                 :::1:39263             ESTABLISHED
-
tcp6       0      0 :::1:39263              :::1:22                ESTABLISHED
3570/sah
```



本章内容

- ☐ 统计磁盘使用情况
- ☐ 计算命令的执行时间
- ☐ 登录用户、启动日志及启动故障的相关信息
- ☐ 打印出10条最常使用的命令
- ☐ 列出1小时内占用CPU最多的10个进程
- ☐ 用watch监视命令输出
- ☐ 对文件及目录访问进行记录
- ☐ 用logrotate管理日志文件
- ☐ 用syslog记录日志
- ☐ 通过监视用户登录找出入侵者
- ☐ 监视远程磁盘的健康情况
- ☐ 找出系统中用户的活动时段

8.1 简介

操作系统是由一系列用于不同目的、服务于不同任务的系统软件组成的。为了了解这些软件是否工作正常，需要由操作系统或系统管理员对其进行监视。我们还使用了一项被称为日志记录(logging)的技术，借助这项技术，应用程序在运行的时候，重要的信息会被写入文件中。通过阅读这个文件，我们就能够得知特定软件或守护进程的操作过程。如果应用程序或服务发生了崩溃，这些信息有助于我们进行调试和故障修复。日志记录和监视也能够帮助我们大量数据中收集信息。为了确保操作系统安全以及排除故障，它们可是相当重要。

在本章中，我们要和监视系统活动的各种命令打交道，同时还要学习日志技术及其使用方法。

8.2 统计磁盘的使用情况

磁盘空间是一种有限资源。我们经常要统计硬盘或其他存储介质的使用情况，以便确定磁盘上的可用空间。当磁盘可用空间开始变少时，我们就得找到大体积的文件删除或移走它们，以便腾出空间。磁盘使用管理多用于shell脚本环境。这则攻略将讲解用于磁盘管理和问题排除的各种命令，通过这些命令的不同选项可以统计磁盘的使用情况。

8.2.1 新手上路

df和du是Linux中用于统计磁盘使用情况的重要命令。df是disk free的缩写，du是disk usage的缩写。让我们看看如何用它们执行各种涉及磁盘使用情况的任务。

8.2.2 实战演练

找出某个文件（或多个文件）占用的磁盘空间：

```
$ du FILENAME1 FILENAME2 ..
```

例如：

```
$ du file.txt
4
```



统计结果默认以字节作为计算单位。

要获得某个目录中所有文件的磁盘使用情况，并在每一行中显示各个文件的磁盘占用详情，可以使用：

```
$ du -a DIRECTORY
```

-a递归地输出指定目录或多个目录中所有文件的统计结果。



执行du DIRECTORY也可以输出类似的结果，但是它只会显示子目录使用的磁盘空间，而不显示每个文件占用空间的情况。要想显示文件的磁盘使用情况，必须使用 -a。

例如：

```
$ du -a test
4 test/output.txt
4 test/process_log.sh
4 test/pcpu.sh
16 test
```

使用du DIRECTORY的输出如下：

```
$ du test
16 test
```

8.2.3 补充内容

让我们看看du命令的其他用法。

1. 以KB、MB或块（block）为单位显示磁盘使用情况

命令du默认显示文件占用的总字节数，但是以标准单位KB、MB或GB显示磁盘使用情况更方便人们阅读。要采用这种更友好的格式进行打印，可以使用选项 -h：

```
du -h FILENAME
```

例如：

```
$ du -sh test/pcpu.sh
4.0K test/pcpu.sh
# 可以接受多个文件参数
```

或者

```
# du -h DIRECTORY
$ du -h hack/
16K hack/
```

2. 显示磁盘使用总计

假如我们需要计算所有文件或目录总共占用了多少磁盘空间，那么显示单个文件的使用情况显然没什么用。du的选项 -c 可以输出作为命令参数的所有文件和目录的磁盘使用情况总计，它会在输出结果末尾加上一行总计。其使用语法如下：

```
$ du -c FILENAME1 FILENAME2..
```

例如：

```
du -c process_log.sh pcpu.sh
4 process_log.sh
4 pcpu.sh
8 total
```

或者

```
$ du -c DIRECTORY
```

例如

```
$ du -c test/
16 test/
16 total
```

或者

```
$ du -c *.txt
# 通配符
```

-c 可以同 -a、-h 等选项组合使用。如果不使用 -c，也可以得到同样的输出。唯一不同的是选项 -c 会添加一行磁盘使用情况总计。

另一个选项 -s (summarize, 合计) 则只输出合计数据。它可以配合 -h 打印出人们易读的格式。这个命令在实践中经常使用，其语法如下：

```
$ du -s FILES(s)
$ du -sh DIRECTORY
```

例如：

```
$ du -sh slynux
680K slynux
```

3. 用特定的单位打印文件

我们可以强制du使用特定的单位打印磁盘使用情况。举例如下。

□ 打印以字节（默认输出）为单位的文件大小：

```
$ du -b FILE(s)
```

□ 打印以KB为单位的文件大小：

```
$ du -k FILE(s)
```

□ 打印以MB为单位的文件大小：

```
$ du -m FILE(s)
```

□ 打印以指定块为单位的文件大小：

```
$ du -B BLOCK_SIZE FILE(s)
```

其中BLOCK_SIZE以字节为单位。

下面的例子中综合了这些命令：

```
$ du pcpu.sh
4 pcpu.sh
$ du -b pcpu.sh
439 pcpu.sh
$ du -k pcpu.sh
4 pcpu.sh
$ du -m pcpu.sh
1 pcpu.sh
$ du -B 4 pcpu.sh
1024 pcpu.sh
```

4. 从磁盘使用统计中排除部分文件

有时候我们需要从磁盘使用统计中排除部分文件。可以使用两种方法指定需要排除的文件：

(1) 通配符

按照下面的方法指定一个通配符：

```
$ du --exclude "WILDCARD" DIRECTORY
```

例如：

```
$ du --exclude "*.txt" FILES(s)
# 排除所有的.txt文件
```

(2) 排除列表

从文件中获取需要排除的文件列表：

```
$ du --exclude-from EXCLUDE.txt DIRECTORY
# EXCLUDE.txt包含了需要排除的文件列表
```

还有其他的选项可以很方便地用来限制磁盘使用统计的范围。我们可以用 `--max-depth` 指定du应该遍历的目录层次的最大深度。将深度指定为1，可以统计当前目录下的所有文件占用内存的情况。将深度指定为2，可以统计当前目录以及下一级子目录中文件占用内存的情况，但不包括第二级子目录。

例如：

```
$ du --max-depth 2 DIRECTORY
```



可以用选项 `-x` 来限制 `du`，使它只能够对单个文件系统进行遍历。假设运行 `du DIRECTORY`，它会递归地遍历每一个可能存在的子目录。目录层次中的某个子目录可能就是一个挂载点（例如，`/mnt/sdai` 是 `/mnt` 的子目录，同时也是设备 `/dev/sdai` 的挂载点）。`du` 会遍历挂载点并统计设备文件系统的磁盘使用情况。为了避免 `du` 的这种行为，可以在使用 `du` 其他选项的同时加上选项 `-x`。`du -x /` 会将 `/mnt` 中的所有挂载点排除在磁盘使用统计之外。

当使用 `du` 时，要确保对其遍历的目录和文件拥有适合的读权限。

5. 找出指定目录中最大的10个文件

我们经常会碰上找出大体积文件的活儿，一般是要将这些大文件删除或移走。我们可以用 `du` 和 `sort` 轻松地完成这项任务。下面的单行脚本就是其实现方法：

```
$ du -ak SOURCE_DIR | sort -nrk 1 | head
```

其中，`-a` 指定了所有的目录和文件。因此 `du` 会遍历 `SOURCE_DIR` 并计算所有文件的大小。由于指定了选项 `-k`，输出的第一列会包含以 KB 为单位的文件大小，第二列则会包含文件或文件夹的名称。`sort` 对第一列依数值逆序排序。`head` 用来显示前 10 行输出。

例如：

```
$ du -ak /home/slynux | sort -nrk 1 | head -n 4
50220 /home/slynux
43296 /home/slynux/.mozilla
43284 /home/slynux/.mozilla/firefox
43276 /home/slynux/.mozilla/firefox/8c22khxc.default
```

上面这个单行脚本的缺点在于它在结果中包含了目录。如果我们只是需要找出最大的文件而不是目录的话，我们可以按照下面的方法改进脚本，使它只输出最大的文件：

```
$ find . -type f -exec du -k {} \; | sort -nrk 1 | head
```

我们利用 `find` 替 `du` 将文件过滤出来，而无需 `du` 本身进行递归遍历。

6. 磁盘可用空间信息

`du` 提供磁盘使用情况信息，而 `df` 提供磁盘可用空间信息。该命令用或不用选项 `-h` 皆可。如果使用 `-h`，`df` 则会以易读的格式打印磁盘空间信息。

例如：

```
$ df
Filesystem      1K-blocks    Used Available Use% Mounted on
/dev/sdai        9611492   2276840   6846412   25% /
none             508828      240    508588    1% /dev
none            513048      168    512880    1% /dev/shm
none            513048       88    512960    1% /var/run
none            513048        0    513048    0% /var/lock
```

```

none          513048      0      513048      0% /lib/init/rw
none          9611492    2276840    6846412    25% /var/lib/
ureadahead/debugfs

$ df -h
FilesystemSize  Used    Avail    Use%    Mounted on
/dev/sdal      9.2G    2.2G    6.6G    25% /
none          497M    240K    497M    1% /dev
none          502M    168K    501M    1% /dev/shm
none          502M    88K     501M    1% /var/run
none          502M    0       502M    0% /var/lock
none          502M    0       502M    0% /lib/init/rw
none          9.2G    2.2G    6.6G    25% /var/lib/ureadahead/debugfs

```

8.3 计算命令执行时间

当测试一个应用程序或比较不同的算法时，程序的执行时间非常重要。一个好算法的执行时间应该最短。在某些情况下，我们需要监视程序执行所需要的时间。例如，当学习排序算法时，你该怎么样表明哪一种算法更快一些？答案就是计算这些算法针对同一个数据集所耗费的时间。让我们看看如何计算。

实战演练

所有的类UNIX操作系统都包含time命令。你可以将time放在需要计算执行时间的命令之前，例如：

```
$ time COMMAND
```

命令COMMAND会执行并生成输出。除此之外，time命令会将命令的执行时间添加到stderr中。例如：

```

$ time ls
test.txt
next.txt
real    0m0.008s
user    0m0.001s
sys     0m0.003s

```

输出中分别显示了执行该命令所花费的real时间、user时间以及system时间。

- ❑ **real**时间指的是挂钟时间（wall clock time），也就是命令从开始执行到结束的时间。这段时间包括其他进程所占用的时间片（time slice）以及进程被阻塞时所花费的时间（例如，为等待I/O操作完成所用的时间）。
- ❑ **user**时间是指进程花费在用户模式（user-mode）中的CPU时间。这是唯一真正用于执行进程所花费的时间。执行其他进程以及花费在阻塞状态中的时间并没有计算在内。
- ❑ **sys**时间是指进程花费在内核模式（in the kernel）中的CPU时间。它代表在内核中执行系统调用所使用的时间，这和库代码（library code）不同，后者仍旧运行在用户空间（user space）。与“user时间”类似，这也是真正由进程使用的CPU时间。



time命令的可执行二进制文件位于/usr/bin/time，还有一个Shell内建命令也叫time。当运行time时，默认调用的是Shell的内建命令。Shell内建的time命令选项有限。因此，如果我们需要使用另外的功能，就应该使用可执行文件time的绝对路径(/usr/bin/time)。

可以用选项-o filename将命令执行时间写入文件：

```
$ /usr/bin/time -o output.txt COMMAND
```

文件名应该出现在选项-o之后。

要将命令执行时间添加到文件而不影响其原有内容，使用选项-a以及-o：

```
$ /usr/bin/time -a -o output.txt COMMAND
```

我们也可以使用选项-f，利用格式字符串格式化时间输出。格式字符串由对应于特定选项的参数组成，这些参数以%作为前缀。real时间、user时间、sys时间的格式字符串分别如下：

☐ real -%e

☐ user -%U

☐ sys -%S

通过结合参数字符串，我们就可以创建格式化输出：

```
$ /usr/bin/time -f "FORMAT STRING" COMMAND
```

例如：

```
$ /usr/bin/time -f "Time: %U" -a -o timing.log uname
Linux
```

其中，%U是user时间的参数。

格式化输出生成后被写入标准输出，执行时间信息被写入标准错误。我们可以用重定向操作符(>)对格式化输出重定向，用错误重定向操作符(2>)对时间信息重定向。例如：

```
$ /usr/bin/time -f "Time: %U" uname > command_output.txt 2>time.log
$ cat time.log
Time: 0.00
$ cat command_output.txt
Linux
```

通过time命令可以获得进程的很多细节信息，包括退出状态、接收到的信号数量、进程上下文的切换次数等。每一个参数都可以用适合的格式字符串显示。

表8-1展示了一些可以使用的参数。

表 8-1

参 数	描 述
%C	进行计时的命令名称以及命令行参数
%D	进程非共享数据区域的大小，以KB为单位
%E	进程使用的real时间（挂钟时间），显示格式为[小时:]分钟:秒
%x	命令的退出状态

(续)

参 数	描 述
%k	进程接收到的信号数量
%W	进程被交换出主存的次数
%Z	系统的页面大小。这是一个系统常量,但在不同的系统中,这个常量值也不同
%P	进程所获得的CPU时间百分比。这个值等于user+system时间除以总的运行时间。结果以百分比形式显示
%K	进程的平均总 (data+stack+text) 内存使用量,以KB为单位
%w	进程主动 (voluntarily) 进行上下文切换的次数,例如等待I/O操作完成
%c	进程被迫 (involuntarily) 进行上下文切换的次数 (由于时间片到期)

例如,用参数 %Z 显示系统页面大小:

```
$ /usr/bin/time -f "Page size: %Z bytes" ls /dev/null
Page size: 4096 bytes
```

这里并不需要被计时命令的输出 (即ls), 因此为了使其不显示在终端中, 我们将标准输出重定向到了/dev/null设备。

还有很多可用格式字符串, 请阅读man time, 以便获取更多的细节信息。

8.4 与当前登录用户、启动日志及启动故障的相关信息

收集与操作环境、当前登录用户、主机加电时间、启动故障等相关的信息很有用处。这则攻略要考察一些用于收集活动主机信息的命令。

8.4.1 新手上路

这则攻略将介绍以下命令: who、w、users、uptime、last和lastb。

8.4.2 实战演练

要获取当前登录用户的相关信息, 可以使用:

```
$ who
slynux pts/0 2010-09-29 05:24 (slynux-macbook-pro.local)
slynux tty7 2010-09-29 07:08 (:0)
```

或者

```
$ w
07:09:05 up 1:45, 2 users, load average: 0.12, 0.06, 0.02
USER TTY FROM LOGIN IDLE JCPU PCPU WHAT
slynux pts/0 slynux 05:24 0.00s 0.65s 0.11s sshd: slynux
slynux tty7 :0 07:08 1:45m 3.28s 0.26s gnome-session
```

该命令会提供当前登录的用户、用户所使用的伪终端 (pseudo TTY)、伪终端当前所执行的命令以及用户登录的IP地址。如果是本地主机, 还会显示主机名。who和w的命令输入格式略有不同。w提供了更多细节。

TTY是与文本终端相关联的设备文件。当用户生成一个新终端时, 对应的设备文件就会出现在/dev/中 (例如/dev/pts/3)。可以通过输入并执行命令tty来获得当前终端的设备路径。

要列出当前登录主机的用户列表, 可以使用:

```
$ users
slynux slynux slynux hacker
```

如果用户打开了多个伪终端, 那么同一个用户会多次显示。在上面的输出中, 用户slynux打开了3个伪终端。排除重复用户出现的最简单的方法是使用sort和uniq进行过滤:

```
$ users | tr ' ' '\n' | sort | uniq
slynux
hacker
```

利用tr将' '替换成'\n', 然后用sort和uniq为每个用户生成唯一的输出。

要查看系统已经通电运行了多长时间, 可以使用:

```
$ uptime
21:44:33 up 3:17, 8 users, load average: 0.09, 0.14, 0.09
```

单词up之后的时间表明了系统已经通电运行了多久。我们可以编写一个简单的单行脚本来提取运行时间。

uptime输出中的平均负载 (load average) 是表明系统负载量的一个参数。在第9章我们会对此详细地进行解释。要获取前一次的启动及用户登录会话的信息, 可以使用:

```
$ last
slynux    tty7          :0                Tue Sep 28 18:27    still logged in
reboot    system boot      2.6.32-21-generi  Tue Sep 28 18:10 - 21:46 (03:35)
slynux    pts/0           :0.0             Tue Sep 28 05:31 - crash (12:39)
```

last命令可以提供登录会话信息。它实际上是一个系统登录日志, 包括了登录tty、登录时间、状态等信息。

last命令以日志文件/var/log/wtmp作为输入日志数据。它也可以用选项-f明确地指定日志文件。例如:

```
$ last -f /var/log/wtmp
```

要获取单个用户登录会话的信息, 可以使用:

```
$ last USER
```

获取重启会话信息:

```
$ last reboot
reboot    system boot 2.6.32-21-generi Tue Sep 28 18:10 - 21:48 (03:37)
reboot    system boot 2.6.32-21-generi Tue Sep 28 05:14 - 21:48 (16:33)
```

获取失败的用户登录会话信息:

```
# lastb
test tty8 :0 Wed Dec 15 03:56 - 03:56 (00:00)
slymox tty8 :0 Wed Dec 15 03:55 - 03:55 (00:00)
```

你必须以超级用户的身份运行lastb。

8.5 打印出 10 条最常使用的命令

终端是用来访问shell的工具，在shell中我们可以输入并执行命令。用户要在shell中运行很多命令，其中某些命令使用得比较频繁。用户的使用习惯可以通过查看其经常使用的命令轻松地得以确认。这则攻略是一个小练习：找出10条最常使用的命令。

8.5.1 新手上路

Bash跟踪用户之前输入过的命令，并将其存储在文件 ~/.bash_history中。不过，它只保留特定数量（例如500）的最近执行过的命令。可以用命令histroy或cat ~/.bash_history查看命令历史记录。我们就用这种方法来找出频繁使用的命令。

8.5.2 实战演练

我们可以从 ~/.bash_history中获得命令列表，只需要命令名称即可，而不需要命令参数，接着统计每条命令的出现次数，然后找出出现次数最多的10条命令。

实现脚本如下：

```
#!/bin/bash
#文件名: top10_commands.sh
#用途: 列出最常使用的10条命令

printf "COMMAND\tCOUNT\n" ;

cat ~/.bash_history | awk '{ list[$1]++; } \
END{
for(i in list)
{
printf("%s\t%d\n",i,list[i]); }
}' | sort -nrk 2 | head
```

输出样例：

```
$ ./top10_commands.sh
COMMAND COUNT
ping      80
ls         56
cat        35
ps         34
sudo       26
du          26
cd          26
ssh         22
sftp        22
clear       21
```



8.5.3 工作原理

在上面的脚本中，用 `~/bash_history` 作为源文件。通过管道将源输入传递给 `awk`。在 `awk` 中，我们使用了一个关联数组。这个数组将命令名作为索引，将命令出现的次数作为数组元素值。命令每出现一次，计数值增加 1 (`list[$1]++`)。\$1 是输入文本行的第一个单词。如果使用 \$0，则包含输入文本行中所有的单词^①。比如说，如果 `ssh 192.168.0.4` 是 `bash_history` 其中的一行，那么 \$0 中包含 `ssh 192.168.0.4`，\$1 中包含 `ssh`。

一旦对命令历史文件中的所有行遍历完毕，我们就得到了一个数组，其中数组索引是命令名称，数组元素值是命令出现的次数。因此出现次数最多的命令名称就是使用最频繁的命令。在 `awk` 的 `END{ }` 语句块中，我们遍历关联数组，打印出所有的命令名及其对应的出现次数。`sort-nrk 2` 对第 2 列 (COUNT) 按照数值逆序排序。最后用 `head` 命令从列表中提取出前 10 条命令。你可以用 `head -n NUMBER` 将前 10 条命令改为前 5 条或是其他任意数量。

8.6 列出 1 小时内占用 CPU 最多的 10 个进程

CPU 时间是一项重要的资源，有时我们需要跟踪某个时期内占用 CPU 周期最多的进程。在普通的桌面系统或膝上系统中，CPU 处于高负荷状态也许不会引发问题。但对于需要处理大量请求的服务器来说，CPU 是极其重要的资源。通过监视某个时期内 CPU 的使用情况，我们可以找出长期占用 CPU 的进程并对其进行优化，或是调试其他问题。这则攻略是一次有关进程监视与日志记录的实践。

8.6.1 新手上路

`ps` 命令用于收集系统中进程的详细信息。这些信息包括 CPU 使用情况、正在执行的命令、内存使用、进程状态等。记录在一小时内占用过 CPU 的进程，然后通过恰当地运用 `ps` 以及文本处理就可以找出占用 CPU 最多的 10 个进程。关于 `ps` 命令的更多细节，请参考第 9 章。

8.6.2 实战演练

让我们看看用于监视并计算一小时内 CPU 使用情况的 shell 脚本：

```
#!/bin/bash
# 文件名: pcpu_usage.sh
# 用途: 计算 1 个小时内进程的 CPU 占用情况

SECS=3600
UNIT_TIME=60

# 将 SECS 更改成需要进行监视的总秒数
```

^① 在 `shell` 和 `awk` 中，都有一个特殊变量 \$0。尽管变量名相同，但是含义却不同。对于 `shell` 来说，\$0 中包含的是命令行的第一个单词，也就是命令名。对于 `awk` 来说，\$0 中包含当前记录。

```

#UNIT_TIME是取样的时间间隔,单位是秒

STEPS=$(( $SECS / $UNIT_TIME ))

echo Watching CPU usage... ;

for((i=0;i<STEPS;i++))
do
    ps -eo comm,pcpu | tail -n +2 >> /tmp/cpu_usage.$$
    sleep $UNIT_TIME
done

echo
echo CPU eaters :

cat /tmp/cpu_usage.$$ | \
awk '
{ process[$1]+=$2; }
END{
    for(i in process)
    {
        printf("%-20s %s",i, process[i] ;
    }

}' | sort -nrk 2 | head
rm /tmp/cpu_usage.$$
#删除临时日志文件

```

输出如下:

```

$ ./pcpu_usage.sh
Watching CPU usage...
CPU eaters :
Xorg          20
firefox-bin   15
bash          3
evince        2
pulseaudio    1.0
pcpu.sh       0.3
wpa_supplicant 0
wnck-applet   0
watchdog/0    0
usb-storage   0

```

8.6.3 工作原理

在上面的脚本中,主要的输入源是`ps -eo comm,pcpu`,其中`comm`表示命令名(command name),`pcpu`表示CPU使用率(CUP usage in percent)。该命令输出所有进程名及CPU使用率。每个进程对应一行输出。因为需要监视一小时内CPU的使用情况,所以我们得在一个每次迭代时间为60秒的for循环中不停地用`ps -eo comm,pcpu | tail -n +2`来获取CPU的使用统计数据,并将这些数据添加到文件`/tmp/cpu_usage.$$`中。60秒的迭代时间通过`sleep 60`来提供。这就使得每分钟执行一次ps。

`tail -n +2`用来将ps 输出中的头部和COMMAND %CPU剔除。

cpu_usage.\$\$ 中的 \$\$ 表示当前脚本的进程 ID。假设进程 ID 为 1345，那么在脚本执行时它会被替换成 /tmp/cpu_usage.1345。因为这是一个临时文件，所以我们把它放在 /tmp 中。

统计文件在 1 小时后就准备妥当了，文件中包含了 60 项，分别对应每分钟的进程状态。然后用 awk 求出每个进程总的 CPU 使用情况。我们用了个关联数组来统计 CPU 使用情况。其中进程名作为数组索引。最后根据总的 CPU 使用情况依数值逆序排序，并通过 head 获得前 10 项。

8.6.4 参考

- 4.7 节讲解了 awk 命令。
- 3.14 节讲解了 tail 命令。

8.7 用 watch 监视命令输出

我们可能需要在某段时期内以固定的间隔时间不断监视某个命令的输出。例如在复制大文件时，我们需要看到不断增长的文件的大小。为了做到这一点，新手们一般会重复输入命令并按回车。其实我们可以利用 watch 命令不断地查看输出。这则攻略就讲解了如何通过 watch 实现这种功能。

8.7.1 实战演练

watch 命令可以用来在终端中以固定的间隔监视命令输出。该命令语法如下：

```
$ watch COMMAND
```

例如：

```
$ watch ls
```

或者

```
$ watch 'COMMANDS'
```

例如：

```
$ watch 'ls -l | grep "^d"'
#只列出目录
```

命令默认每 2 秒更新一次输出。

我们可以用 -n SECOND 指定需要更新输出的时间间隔。例如：

```
$ watch -n 5 'ls -l'
#以5秒为间隔，监视ls -l的输出
```

8.7.2 补充内容

让我们研究一下 watch 命令的其他特性。

突出 (highlighting) watch 输出中的差异

watch 有一个选项可以将时间间隔前后的命令输出差异以不同颜色突出标示出来。选项 -d

可以启用这一功能：

```
$ watch -d 'COMMANDS'
```

8.8 对文件及目录访问进行记录

记录文件及目录访问对于跟踪文件和目录的变化很有帮助。这则攻略将讲解如何记录用户访问。

8.8.1 新手上路

inotifywait命令可以用来收集有关文件访问的信息。Linux发行版并没有默认包含这个命令，你得用软件包管理器自行安装inotify-tools。这个命令还需要将inotify支持编译入Linux内核，好在大多数新的GNU/Linux发行版都在内核中启用了inotify。

8.8.2 实战演练

来看看用来监视目录访问的shell脚本：

```
#!/bin/bash
#文件名: watchdir.sh
#用途: 监视目录访问
path=$1
#将目录或文件路径作为脚本参数
inotifywait -m -r -e create,move,delete $path -q
```

输出样例如下：

```
$ ./watchdir.sh .
./ CREATE new
./ MOVED_FROM new
./ MOVED_TO news
./ DELETE news
```

8.8.3 工作原理

上面的脚本记录给定了路径中文件或目录的创建、移动以及删除。选项 `-m` 表明要持续监视变化，而不是在事件发生之后退出。`-r` 允许采用递归形式监视目录。`-e` 指定需要监视的事件列表。`-q` 用于减少冗余信息，只打印出所需要的信息。命令输出可以被重定向到日志文件。

我们可以从事件列表中添加或删除事件。一些重要的事件如表8-2所示。

表 8-2

事 件	描 述
访问 (access)	读取文件
修改 (modify)	文件内容被修改
属性 (attrib)	文件元数据被修改

(续)

事 件	描 述
移动 (move)	对文件进行移动操作
创建 (create)	生成新文件
打开 (open)	对文件进行打开操作
关闭 (close)	对文件进行关闭操作
删除 (delete)	文件被删除

8.9 用 logrotate 管理日志文件

日志文件是Linux系统维护中必不可少的组成部分。日志文件可以帮助跟踪系统中多种服务所发生的事件。这有助于系统管理员排查问题,同时也为活动主机上出现的事件提供了统计数据。随着时间的推移,日志文件会变得越来越庞大,因而对于日志文件的管理必不可少。我们将用一种被称为轮替 (rotation) 的技术来限制日志文件的体积,一旦它超过了限定的大小,就要对它的内容进行抽取 (strip), 同时将日志文件中的旧条目存储到归档文件中。因此旧的日志文件就得以保存以便以后参阅。让我们看看如何对日志文件进行轮替及存储操作。

8.9.1 新手上路

logrotate是每一位Linux系统管理员都应该了解的命令。它能够将日志文件的大小限制在给定的SIZE内。各种应用程序会将信息添加到日志文件中。最新添加的信息总是出现在日志文件的尾部。logrotate根据配置文件扫描特定的日志文件。它在保留日志文件中最新添加的100KB内容 (例如指定SIZE = 100k) 的同时, 将剩下的数据 (较旧的日志数据) 移入新文件logfile_name.1。如果该日志文件 (logfile_name.1) 中的内容越来越多, 已经超出了SIZE的定额, logrotate就会再用最近的内容更新日志文件, 然后用较旧的内容创建logfile_name.2。这个过程可以轻松地使用logrotate进行配置。logrotate还可以将旧的日志文件压缩成logfile_name.1.gz、logfile_name.2.gz等。是否选择压缩旧日志文件也可以通过logrotate来配置。

8.9.2 实战演练

logrotate的配置目录位于/etc/logrotate.d。如果列出这个目录中的内容, 你会发现很多其他的日志文件配置。

我们可以为自己的日志文件编写一个特定的配置:

```
$ cat /etc/logrotate.d/program
/var/log/program.log {
    missingok
    notifempty
    size 30k
    compress
    weekly
    rotate 5
```

```
create 0600 root root
}
```

这就是全部的配置。其中，`/var/log/program.log`指定了日志文件路径。旧的日志文件归档之后也放入同一个目录中。让我们看看这些参数，如表8-3所示。

表 8-3

参 数	描 述
missingok	如果日志文件丢失，则忽略；然后返回（不对日志文件进行轮替）
notifempty	仅当源日志文件非空时才对其进行轮替
size 30k	限制实施轮替的日志文件的大小。可以用1M表示1MB
compress	允许用gzip对较旧的日志进行压缩
weekly	指定进行轮替的时间间隔。可以是weekly、yearly或daily
rotate 5	这是需要保留的旧日志文件的归档数量。在这里指定的是5，所以这些文件名将会是program.log.1.gz、program.log.2.gz等直到program.log.5.gz
create 0600 root root	指定所要创建的归档文件的模式、用户以及用户组

表8-3中的选项都是可选的，我们只在logrotate配置文件中指定所需的选项即可。如果想了解更多的可用选项，请参考logrotate的手册页（<http://linux.die.net/man/8/logrotate>）。

8.10 用 syslog 记录日志

日志文件是那些为用户提供服务的应用程序的重要组成部分。应用程序在运行时将状态信息写入日志文件。如果程序崩溃或者当我们需要查询服务的相关信息时，就可以借助日志文件。在`/var/log`目录中，你可以找到与各种守护进程和应用程序相关的日志文件。`/var/log`是存储日志文件的公共目录。如果你读过日志文件，你就会看出它们都采用了一种通用的格式。在Linux系统中，在`/var/log`中创建并写入日志信息是由被称为syslog的协议处理的。它由守护进程syslogd负责执行。每一个标准应用进程都可以用syslog记录日志信息。在这则攻略中，我们将讨论如何在脚本中用syslogd记录日志信息。

8.10.1 新手上路

日志文件有助于我们推断系统出现了什么故障。因此在编写重要的应用程序时，应当将应用进程的执行过程记录在日志文件中，要养成这个良好的实践习惯。我们接下来将要学习用命令logger通过syslogd记录日志。在学习如何往日志文件中写入信息之前，让我们先看看Linux中一些重要的日志文件，如表8-4所示。

表 8-4

日志文件	描 述
<code>/var/log/boot.log</code>	系统启动信息
<code>/var/log/httpd</code>	Apache Web服务器日志

(续)

日志文件	描 述
/var/log/message	发布内核启动信息
/var/log/auth.log	用户认证日志
/var/log/dmesg	系统启动信息
/var/log/mail.log	邮件服务器日志
/var/log/Xorg.0.log	X服务器日志

8.10.2 实战演练

要向syslog文件 /var/log/message中记录日志信息，可以使用：

```
$ logger LOG_MESSAGE
```

例如：

```
$ logger This is a test log line
```

```
$ tail -n 1 /var/log/messages
```

```
Sep 29 07:47:44 slynux-laptop slynux: This is a test log line
```

/var/log/messages是一个一般用途的日志文件。如果使用logger命令，它默认记录日志信息到 /var/log/messages中。如果要记录特定的标记 (tag)，可以使用：

```
$ logger -t TAG This is a message
```

```
$ tail -n 1 /var/log/messages
```

```
Sep 29 07:48:42 slynux-laptop TAG: This is a message
```

Syslog处理/var/log下的多个日志文件。但是当logger发送消息时，它用标记字符串来确定应该记录到哪一个日志文件中。syslogd使用与日志相关联的TAG来决定应该将其记录到哪一个文件中。你可以从位于/etc/rsyslog.d/目录的配置文件中看到与日志文件相关联的标记字符串。

要将另一个日志文件的最后一行记录到系统日志中，可以使用：

```
$ logger -f /var/log/source.log
```

8.10.3 参考

3.14节讲解了head和tail命令。

8.11 通过监视用户登录找出入侵者

日志文件可以用于收集有关系统状态的细节信息。下面是一个有趣的脚本编写练习。

我们有一个通过SSH连接到Internet的系统。很多攻击者试图登入这个系统，因此需要编写shell脚本来设计一个入侵检测系统。入侵者被定义为：屡次试图登入系统达两分钟以上，并且期间的登录过程全部失败。凡是这类用户都应该被检测出来并生成包含以下细节信息的报告。

□ 试图登录的账户

- 试图登录的次数
- 攻击者的IP地址
- IP地址所对应的主机
- 进行登录过程的时间段

8.11.1 新手上路

我们可以编写这样一个shell脚本，使它能够扫描日志文件并从中收集所需要的信息。我们要处理SSH登录失败的情况。用户认证会话日志被记录在日志文件 `/var/log/auth.log` 中。脚本应该扫描这个日志文件来检测出失败的登录信息，并执行不同的检查来获取所需要的数据。我们可以用 `host` 命令找出IP地址所对应的主机。

8.11.2 实战演练

让我们编写一个入侵检测脚本，它利用用户认证日志文件来生成有关入侵者的报告：

```
#!/bin/bash
#文件名: intruder_detect.sh
#用途: 入侵报告工具, 以auth.log作为日志文件
AUTHLOG=/var/log/auth.log

if [[ -n $1 ]];
then
    AUTHLOG=$1
    echo Using Log file : $AUTHLOG
fi
LOG=/tmp/valid.$$log
grep -v "invalid" $AUTHLOG > $LOG
users=$(grep "Failed password" $LOG | awk '{ print $(NF-5) }' | sort |
uniq)

printf "%-5s|%-10s|%-10s|%-13s|%-33s|%-33s\n" "Sr#" "User" "Attempts" "IP
address" "Host_Mapping" "Time range"

ucount=0;
ip_list=$(egrep -o "[0-9]+\.[0-9]+\.[0-9]+\.[0-9]+" $LOG | sort |
uniq)

for ip in $ip_list;
do
    grep $ip $LOG > /tmp/temp.$$log
for user in $users;
do
    grep $user /tmp/temp.$$log > /tmp/$$.log
    cut -c-16 /tmp/$$.log > $$time
    tstart=$(head -1 $$time);
    start=$(date -d "$tstart" "+%s");
    tend=$(tail -1 $$time);
    end=$(date -d "$tend" "+%s");
    limit=$(( $end - $start ))
    if [ $limit -gt 120 ];
```

```

then
    let ucount++;

    IP=$(egrep -o "[0-9]+\.[0-9]+\.[0-9]+\.[0-9]+" /tmp/$$.log | head
-1 );
    TIME_RANGE="$tstart-->$tend"
    ATTEMPTS=$(cat /tmp/$$.log|wc -l);
    HOST=$(host $IP | awk '{ print $NF }' )

    printf "%-5s|%-10s|%-10s|%-10s|%-33s|%-s\n" "$ucount" "$user"
"$ATTEMPTS" "$IP" "$HOST" "$TIME_RANGE";
    fi
done
done
rm /tmp/valid.$$.log /tmp/$$.log $$time /tmp/temp.$$.log 2> /dev/null

输出样例如下:

```

```

slynux@slynux-laptop:~$ ./intruder_detect.sh sampleauth.log
Using Log file : sampleauth.log

```

Sr#	User	Attempts	IP address	Host Mapping	Time range
1	alice	3	203.110.250.34	attk1.foo.com	Oct 29 05:28:59 -->Oct 29 05:31:59
2	bob1	3	203.110.251.31	attk2.foo.com	Oct 29 05:21:52 -->Oct 29 05:29:52
3	bob2	3	203.110.250.34	attk1.foo.com	Oct 29 05:22:59 -->Oct 29 05:25:52
4	gvvrju	20	203.110.251.31	attk2.foo.com	Oct 28 04:37:10 -->Oct 29 05:19:09
5	root	21	203.110.253.32	attk3.foo.com	Oct 29 05:18:01 -->Oct 29 05:37:81

8.11.3 工作原理

在intruder_detect.sh脚本中，我们将auth.log文件作为输入。我们也可以用脚本的命令行参数来提供一个日志文件作为输入，或者默认读取 /var/log/auth.log。我们只需要记录合法用户的登录日志详情。如果有非法用户企图登录，则类似于“Failed password for invalid user bob from 203.83.248.32 port 7016 ssh2”的日志就会被记录在auth.log中。因此我们需要排除日志文件中所有包含“invalid”的行。grep命令的反转选项（-v）可以用来移除对应于非法用户的所有日志内容。下一步是找出试图登录并且失败的用户列表。对于密码错误，SSH会记录类似的日志信息：

ssh[21197]: Failed password for bob1 from 203.83.248.32 port 50035 ssh2.

所以应找出所有包含“failed password”的行，并找出所有不重复的IP地址以提取对应于每一个IP地址的日志行。提取IP地址列表可以用匹配IP地址的正则表达式和grep命令来完成。用for循环对IP地址进行迭代，并用grep找出对应的日志行并将其写入临时文件。日志行中倒数第6个单词是用户名（例如bob1），可用awk命令提取用户名（倒数第6个单词）。NF返回最后一个单词的列号，因此NF-5就是倒数第6个单词的列号。我们再用sort和uniq生成一个没有重复的用户列表。

现在我们要收集表明登录失败的日志行，这些行中包含了用户名。for循环用来读取对应每一位用户的日志行，并将这些行写入临时文件。每一行的前16个字符是时间戳，可用cut命令提取时间戳。一旦我们得到了一个用户所有登录失败的时间戳，就要检查第一次和最后一次试图登录之间的时间差。第一行日志对应第一次登录，最后一行日志对应最后一次登录。我们用head -1提取首行，用tail -1提取末行。我们也有了首次登录和末次登录的字符串格式的时间戳

(tstart和tends)。使用date命令，我们可以将字符串形式的日期转换成UNIX纪元时的总秒数(1.9节讲解了UNIX纪元时)。

变量start和end中包含着秒数，对应于字符串形式起始时间戳。对这两个时间求差，并检查差值是否大于2分钟（即120秒）。如果某个用户被确认为入侵者，其对应的细节信息就要被生成日志。IP地址可以用正则表达式和egrep命令从日志中提取。试图登录的次数就是某个用户的日志行数。这个行数可以用wc命令获得。IP地址到主机名的映射可以通过将IP地址作为host命令的参数，然后从命令输出中提取。时间范围可以用我们已经提取到的时间戳来打印。最后，删除脚本使用过的临时文件。

上面的脚本旨在演示一个用于扫描日志并从中生成报告的模型。我们尽力让脚本更短小简单，以避免过于复杂，因此难免存在一些bug。你可以使用更好的逻辑来改进脚本。

8.12 监视远程磁盘的健康情况

网络是由不同用户的多台主机组成。它需要对远程主机的磁盘使用情况进行集中监视。网络系统管理员得每天记录网络中所有主机的磁盘使用。日志的每一行应该包含日期、主机IP地址、设备、设备容量、占用空间、剩余空间、使用比例、健康状况等细节信息。如果任何远程主机中的任意分区使用率超过了80%，那么健康状态应该被设置为ALERT，否则就可以设置为SAFT。这则攻略将演示如何编写一个可以收集网络远程主机详细信息的监视脚本。

8.12.1 新手上路

我们需要分别从网络中每台主机收集磁盘使用情况的统计信息，并写入中央主机的日志文件中。可以将负责收集信息并写入日志的脚本调度为每天执行。SSH可用于登录远程系统收集磁盘使用数据。

8.12.2 实战演练

首先，我们得在所有网络中的远程主机上设置一个共用账户。这个账户供脚本disklog登录系统使用。我们需要为这个账户配置SSH自动登录（7.7节讲解了自动登录的配置方法）。我们假设在所有配置了自动登录的远程主机上都有一个叫做test的用户，那么来看看这个shell脚本：

```
#!/bin/bash
# 文件名: disklog.sh
# 用途: 监视远程系统的磁盘使用情况

logfile="diskusage.log"
if [[ -n $1 ]]
then
    logfile=$1
fi
if { ! -e $logfile }
then
```

```

    printf "%-8s %-14s %-9s %-8s %-6s %-6s %s\n" "Date" "IP
address" "Device" "Capacity" "Used" "Free" "Percent" "Status" >
$logfile
fi

IP_LIST="127.0.0.1 0.0.0.0"
#提供远程主机IP地址列表
(
for ip in $IP_LIST;
do
    ssh slynux@$ip 'df -H' | grep ^/dev/ > /tmp/$$.df

    while read line;
    do
        cur_date=$(date +%D)
        printf "%-8s %-14s " $cur_date $ip
        echo $line | awk '{ printf("%-9s %-8s %-6s %-6s
%-8s", $1, $2, $3, $4, $5); }'

        pusg=$(echo $line | egrep -o "[0-9]+%")
        pusg=${pusg/%\%/};
        if [ $pusg -lt 80 ];
        then
            echo SAFE
        else
            echo ALERT
        fi
    done< /tmp/$$.df
done
) >> $logfile

```

我们可以用cron以固定的间隔来调度脚本执行，例如在crontab中写入以下条目，以便实现每天上午10点运行脚本：

```
00 10 * * * /home/path/disklog.sh /home/user/diskusg.log
```

执行命令crontab -e。添加上面一行并保存。

可以手动执行脚本：

```
$ ./disklog.sh
```

脚本的输出样例如下：

Date	IP address	Device	Capacity	Used	Free	Percent	Status
12/15/10	127.0.0.1	/dev/sda1	9.9G	2.4G	7.0G	26%	SAFE
12/15/10	0.0.0.0	/dev/sda1	9.9G	2.4G	7.0G	26%	SAFE

8.12.3 工作原理

在脚本disklog.sh中，我们可以提供日志文件路径作为命令行参数，否则脚本使用默认的日志文件。如果日志文件不存在，它会将日志文件头部写入新文件中。-e logfile用来检查文件是否存在。远程主机的IP地址列表被存储在变量IP_LIST中，彼此之间以空格分隔。要确保在

IP_LIST中列出的所有远程系统中都有用户test, 并且SSH已经配置了自动登录。for循环用来对IP地址进行逐个迭代。通过ssh使用远程命令df -H获取磁盘剩余空间。这项数据被存储在一个临时文件中, while循环用来逐个读取这个文件。用awk提取并打印数据, 同时一并打印的还有日期。用egrep提取使用率, 并将%删除以获取使用率的数值部分。检查得到的数值, 看是否超过了80。如果不足80, 将状态设置为SAFT; 如果大于或等于80, 则将状态设置为ALERT。打印出来的所有数据要被重定向到日志文件中。因此代码被放入子shell()中, 并将标准输出重定向到日志文件。

8.12.4 参考

9.8节讲解了crontab命令。

8.13 找出系统中用户的活动时段

考虑一个使用共享主机(shared hosting)的Web服务器。每天都有很多用户登录和注销。用户活动都被记入服务器的系统日志。这则攻略是一项实践任务: 利用系统日志找出每个用户在服务器上停留了多久, 并根据用户的使用时间对他们进行分级, 最后生成一份包含等级、用户名、首次登录时间、末次登录时间、登录次数以及总时长等细节信息的报告。让我们看看如何解决这个问题。

8.13.1 新手上路

last命令用来列出一个系统中有关用户登录会话的细节。这些会话数据被存储在/var/log/wtmp文件中。通过分别累计各用户的会话时间, 我们就能得出各用户的总使用时间。

8.13.2 实战演练

研究一下这个用来找出活跃用户并生成报告的脚本:

```
#!/bin/bash
#用户名: active_users.sh
#用途: 查找活跃用户

log=/var/log/wtmp
if [[ -n $1 ]];
then
    log=$1
fi
printf "%-4s %-10s %-10s %-6s %-8s\n" "Rank" "User" "Start" "Logins"
"Usage hours"
last -f $log | head -n -2 > /tmp/ulog.$$
cat /tmp/ulog.$$ | cut -d' ' -f1 | sort | uniq> /tmp/users.$$
{
    while read user;
    do
        grep ^$user /tmp/ulog.$$ > /tmp/user.$$
```



```

seconds=0
while read t
do
s=$(date -d $t +%s 2> /dev/null)
let seconds=seconds+s
done< <(cat /tmp/user.$$ | awk '{ print $NF }' | tr -d ' '){}

firstlog=$(tail -n 1 /tmp/user.$$ | awk '{ print $5,$6 }')
nlogins=$(cat /tmp/user.$$ | wc -l)
hours=$(echo "$seconds / 60.0" | bc)
printf "%-10s %-10s %-6s %-8s\n" $user "$firstlog" $nlogins $hours
done< /tmp/users.$$
) | sort -nrk 4 | awk '{ printf("%-4s %s\n", NR, $0) }'
rm /tmp/users.$$ /tmp/user.$$ /tmp/ulog.$$

```

样例输出如下：

```

$ ./active_users.sh
Rank  User      Start   Logins  Usage hours
1     easyibaa  Dec 11   531    11437311943
2     demoproj  Dec 10   350    7538718253
3     kjayaram  Dec 9    213    4587849555
4     cinenews  Dec 11   85     1830831769
5     thebenga  Dec 10   54     1163118745
6     gateway2 Dec 11   52     1120038550
7     soft132  Dec 12   49     1055420578
8     sarathla  Nov 1    45     969268728
9     gtminis  Dec 11   41     883107030
10    agentcde Dec 13   39     840029414

```

8.13.3 工作原理

在脚本active_users.sh中，我们要提供日志文件wtmp作为命令行参数，也可以使用日志文件defaultwtmp。命令last-f用来打印日志文件的内容。日志文件的第一列是用户名。我们用cut从中提取第一列，然后用sort和uniq找出不重复的用户。对每一位用户，用grep找出其对应登录会话的日志行并写入一个临时文件。日志的最后一列是用户登录会话的时长。为了找出用户总的使用时间，需要累加所有的会话时长。使用时间的格式是（小时：秒），因此还需要用date命令将其转换成秒。

要提取用户的会话时长，得使用awk命令。要移除括号，得使用tr -d。用<(COMMANDS)操作符将用户使用时长字符串列表作为标准输入传递给while循环。它的作用就像是文件输入。利用date命令将每一个时长字符串转换成秒数，并累加到变量seconds中。将出现在最后一行的用户的首次登录时间提取出来。登录次数就是日志的行数。要根据总的使用时间来计算每位用户的等级，数据记录需要将总使用时间作为键，进行降序排列。用sort命令的-nr选项指定按照数值逆序排列。-k4指定键的列号（即使用时间）。最后，sort的输出被传递给awk。awk命令为每一行添加上行号，这个行号就是每一位用户的等级。

本章内容

- 收集进程信息
- 杀死进程以及发送或响应信号
- which、whereis、file、whatis与平均负载
- 向用户终端发送消息
- 收集系统信息
- 用/proc收集信息
- 用cron进行调度
- 从Bash中读写MySQL数据库
- 用户管理脚本
- 图像文件的批量缩放及格式转换

9.1 简介

GNU/Linux的生态系统是由运行的程序、服务、连接的设备、文件系统、用户等组成的。按照我们需要的方式对整个系统有一个概观并对操作系统进行整体上的管理，这是系统管理的主要目的。我们应该掌握用于收集系统信息及管理资源的常用命令以及适合的实践用法，以编写执行管理任务的脚本和自动化工具。本章将介绍一些收集系统信息的命令和方法，还会讲解如何利用这些命令编写管理脚本。

9.2 收集进程信息

进程是程序的运行实例（running instance）。运行在一台计算机中的多个进程各自都分配了一个称为进程ID的唯一标识数字。这个数字是一个整数。同一个程序的多个实例可以同时运行，但是他们的进程ID却互不相同。一个进程包括多种属性，例如拥有该进程的用户、进程使用的内存数量、进程占用的CPU等。这则攻略将学习如何收集有关进程的信息。

9.2.1 新手上路

和进程管理相关的重要命令是top、ps和pgrep。让我们看看如何收集进程信息。

9.2.2 实战演练

ps是收集进程信息的重要工具。它提供的信息包括：拥有进程的用户、进程的起始时间、进程所对应的命令行路径、进程ID (PID)、进程所属的终端 (TTY)、进程使用的内存、进程占用的CPU等。例如：

```
$ ps
  PID TTY          TIME CMD
 1220 pts/0    00:00:00 bash
 1242 pts/0    00:00:00 ps
```

ps命令通常结合一系列参数使用。如果不使用任何参数，ps将显示运行在当前终端 (TTY) 中的进程。第一列显示进程ID (PID)，第二列是TTY (终端)，第三列是进程启动后过去的时间，最后一列是CMD (进程所对应的命令)。

为了显示包含更多信息的更多列，可以使用-f，如下所示。

```
$ ps -f
  UID          PID    PPID  C  STIME TTY          TIME CMD
slynnux      1220    1219   0   18:18 pts/0    00:00:00 -bash
slynnux      1587    1220   0   18:59 pts/0    00:00:00 ps -f
```

上面的ps命令没有什么用处，因为它没有提供当前终端外的任何进程信息。要获取运行在系统中的每一个进程的信息，使用选项-e (every)。选项-ax (all) 也可以生成同样的输出。



选项-x和-a指定解除由ps默认添加的TTY限制。通常使用不带参数的ps命令，以便只打印出其所属终端的进程。

运行ps -e或ps -ef，要么是ps -ax或ps -axf；

```
$ ps -e | head
  PID TTY          TIME CMD
 1 ?          00:00:00 init
 2 ?          00:00:00 kthreadd
 3 ?          00:00:00 migration/0
 4 ?          00:00:00 ksoftirqd/0
 5 ?          00:00:00 watchdog/0
 6 ?          00:00:00 events/0
 7 ?          00:00:00 cpuset
 8 ?          00:00:00 khelper
 9 ?          00:00:00 netns
```

输出列表很长。我们使用head进行了过滤，所以只列出了前10项。

ps命令支持显示除进程名及进程ID之外的多种信息。ps默认在不同的列中显示这些信息。这些信息中的大多数对我们来说没什么用处。我们可以用-o来指定想要显示的列。这样就可以只打印出我们需要的内容。与进程相关的参数可以通过与此参数对应的命令选项指定。参数列表以及-o的用法在接下来会进行讨论。

用ps显示需要的输出列：

```
$ ps [OTHER OPTIONS] -o parameter1,parameter2,parameter3 ..
```

-o的参数以逗号操作符(,)作为定界符。值得注意的是逗号操作符与它分隔的参数之间是没有空格的。在大多数情况下,选项-o都是和选项-e (every)结合使用的(-oe),因为它需要列出运行在系统中的每一个进程。但是如果-o需要使用某些过滤器,例如列出特定用户拥有的进程,那么就不再使用-e。-e和过滤器结合使用将没有任何实际效果,依旧会显示所有的进程。

示例如下。其中comm表示COMMAND,pcpu表示CPU占用率:

```
$ ps -eo comm,pcpu | head
COMMAND      %CPU
Init          0.0
kthreadd      0.0
migration/0   0.0
ksoftirqd/0   0.0
watchdog/0    0.0
events/0      0.0
cpuset        0.0
khelper       0.0
netns         0.0
```

选项-o可以使用不同的参数,这些参数及其描述如表9-1所示。

表 9-1

参 数	描 述
pcpu	CPU占用率
pid	进程ID
ppid	父进程ID
pmem	内存使用率
comm	可执行文件名
cmd	简单命令 (simple command) ^①
user	启动进程的用户
nice	优先级 (niceness)
time	累计的CPU时间
etime	进程启动后度过的时间
tty	所关联的TTY设备
euid	有效用户ID
stat	进程状态

① 简单命令是我们平时使用最频繁的一种命令。它是由空白字符分隔的一系列单词,以shell控制操作符作为结尾。第一个单词指定要执行的命令,余下的单词作为命令参数。Shell控制操作符可以是换行符,或者是:|, &&, &, ;, ;, ;, |, &, {, }。

9.2.3 补充内容

让我们看看其他进程控制命令的用例。

1. top

top对于系统管理员来说是一个极为重要的命令。top命令默认会输出一个占用CPU最多的进程列表。该命令的用法如下：

```
$ top
```

除了若干个占用CPU最多的进程外，该命令还会显示一些其他参数。

2. 根据参数对ps输出进行排序

可以用--sort将ps命令的输出根据特定的列进行排序。

可以在参数前加上+（升序）或-（降序）来指定排序方式：

```
$ ps [OPTIONS] --sort -parameter1,+parameter2,parameter3..
```

例如，要列出占用CPU最多的10个进程，可以使用：

```
$ ps -eo comm,pcpu --sort -pcpu | head
COMMAND          %CPU
Xorg              0.1
hald-addon-stor   0.0
ata/0             0.0
scsi_eh_0         0.0
gnome-settings-   0.0
init              0.0
hald              0.0
pulseaudio        0.0
gdm-simple-gree   0.0
```

进程依据CPU占用率进行降序排序，用head命令提取前10个进程。

我们可以用grep从ps的输出中提取与给定进程名或其他参数相关的条目。要找出与bash进程相关的条目，可以使用：

```
$ ps -eo comm,pid,pcpu,pmem | grep bash
bash             1255  0.0  0.3
bash             1680  5.5  0.3
```

3. 找出给定命令名对应的进程ID

假设一个命令有多个实例正在运行，我们可能需要识别这些进程的进程ID。这个信息可以使用ps或pgrep命令得到。按照下面的方式使用ps：

```
$ ps -C COMMAND_NAME
```

或者

```
$ ps -C COMMAND_NAME -o pid=
```

用户自定义格式指示符-o先前已经讲解过了。但是这里你可以看到pid后面加上了=，这将移除ps输出中的头部PID。在参数后加上=就可以移除每一列的头部。例如：

```
$ ps -C bash -o pid=
1255
1680
```

这条命令列出了所有Bash进程的进程ID。

除此之外，还有一个很方便的工具pgrep。你可以用它获得一个特定命令的进程ID列表。例如：

```
$ pgrep COMMAND
$ pgrep bash
1255
1680
```



pgrep只需要命令名的一部分作为输出参数来提取Bash命令，诸如pgrep ash或pgrep bas都能够奏效，但是ps需要你输入命令准确的全名。

pgrep可以接受很多输出过滤选项。如果要指定输出定界符，而不以换行符作为定界符，可以这样：

```
$ pgrep COMMAND -d DELIMITER_STRING
$ pgrep bash -d ":"
1255:1680
```

指定进程的用户（拥有者）列表：

```
$ pgrep -u root,slynux COMMAND
```

其中root和slynux都是用户名。

返回所匹配的进程数量：

```
$ pgrep -c COMMAND
```

4. 根据真实用户或ID以及有效用户或ID过滤ps输出

可以用ps根据指定的真实/有效用户名或ID对进程进行分组。指定的参数可以用来过滤ps的输出：通过检查每一个输出条目是否属于参数列表中指定的有效用户或真实用户，并只显示匹配的条目。实现方法如下：

- ☐ 用 -u EUSER1、EUSER2 依次类推，指定有效用户列表；
- ☐ 用 -U RUSER1、RUSER2 依次类推，指定真实用户列表。

例如：

```
$ ps -u root -U root -o user,pcpu
```

该命令会显示以root作为有效用户ID和真实用户ID的所有进程，以及用户、CPU占用率列。



在大多数情况下，-o都是和-e结合使用的，并写成-eo的形式。但是当使用过滤器的时候，-o应该像上面那样单独使用。

5. 用TTY过滤ps输出

可以通过指定进程所属的TTY选择ps的输出。用选项 `-t` 指定TTY列表：

```
$ ps -t TTY1, TTY2 ..
```

例如：

```
$ ps -t pts/0,pts/1
  PID TTY          TIME CMD
 1238 pts/0      00:00:00 bash
 1835 pts/1      00:00:00 bash
 1864 pts/0      00:00:00 ps
```

6. 进程线程的相关信息

通常，与进程线程相关的信息在ps输出中是看不到的。我们可以用选项 `-L` 在ps输出中显示线程的相关信息。这会显示出两列：NLWP和NLP。NLWP是进程的线程数量，NLP是ps输出中每个条目的线程ID。例如：

```
$ ps -eLf
```

或者

```
$ ps -eLf --sort -nlwp | head
  UID      PID PPID    LWP C   NLWP  STIME TTY          TIME CMD
  root      647   1    647 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    654 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    656 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    657 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    658 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    659 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    660 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    662 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
  root      647   1    663 0     64  14:39 ?           00:00:00 /usr/sbin/
 console-kit-daemon --no-daemon
```

该命令列出了线程数最多的10个进程。

7. 指定输出宽度以及所要显示的列

我们可以使用用户自定义输出格式指示符来指定在ps输出中所要显示的列。另一种指定输出格式的方法是使用“标准”选项。根据你的使用方式进行应用，可以尝试以下选项：

- ☐ `-f` `ps -ef`
- ☐ `u` `ps -e u`
- ☐ `ps` `ps -e w` (w代表宽松输出)

8. 显示进程的环境变量

了解某个进程依赖哪些环境变量，这类重要的信息通常是我们用得着的。无论进程的正常运行是否特别依赖于的一组环境变量。我们都可以利用环境变量调试、修复与进程运行相关的问题。

要伴随ps条目同时列出环境变量，可以使用：

```
$ ps -eo cmd e
```

例如：

```
$ ps -eo pid,cmd e | tail -n 3
1162 hald-addon-acpi: listening on acpid socket /var/run/acpid.socket
1172 sshd: slynux [priv]
1237 sshd: slynux@pts/0
1238 -bash USER=slynux LOGNAME=slynux HOME=/home/slynux PATH=/usr/
local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games
MAIL=/var/mail/slynux SHELL=/bin/bash SSH_CLIENT=10.211.55.2 49277 22
SSH_CONNECTION=10.211.55.2 49277 10.211.55.4 22 SSH_TTY=/dev/pts/0
TERM=xterm-color LANG=en_IN XDG_SESSION_COOKIE=d1e96f5cc8a7a3bc3a0a73e44c
95121a-1286499339.592429-1573657095
```

这种环境跟踪方法的用途之一就是解决apt-get软件包管理器的故障。如果你是用HTTP代理连接Internet，你可能得设置环境变量http_proxy=host:port。但是有时候即使设置了代理服务器，但apt-get却弃之不用，照旧返回一个错误。那么你就需要查看环境变量，跟踪这个问题。

我们可能需要借助crontab这类调度工具来使一些应用程序能够自动运行。但是这些应用也许要依赖某些环境变量。假设我们需要在某个特定时间打开一个窗口化的GUI应用程序。我们用crontab将其调度到指定的时间。然而，如果出现像下面这样的条目，你就会发现这个应用程序并没有按时启动：

```
00 10 * * * /usr/bin/windowapp
```

这是因为窗口化的应用程序总是依赖于环境变量DISPLAY。环境变量需要传递给应用程序。首先手动运行windowapp，然后运行ps -C windowapp -eo cmd e。

找出环境变量，将其添加到crontab中的命令名称之前，这个问题就可以搞定了。

按照下面的方法修改条目：

```
00 10 * * * DISPLAY=:0 /usr/bin/windowapp
```

DISPLAY=:0 可以从ps输出中获取。

9.2.4 参考

9.8节讲解了如何调度任务。

9.3 杀死进程以及发送或响应信号

终结进程是我们通常都会碰到的事儿。有时，我们可能需要终结某个程序的所有实例。命令行提供了多种用于终结程序的方法。在类UNIX环境中与进程相关的一个重要概念就是信号。信

号是一种进程间通信机制，它用来中断运行的进程以执行某些操作。终止程序也是通过使用信号技术来实现的。这则攻略介绍了信号及其用法。

9.3.1 新手上路

信号是Linux中的一种进程间通信机制。我们可以使用特定的信号来中断进程。每一种信号都同一个整数值相关联。当进程接收到一个信号时，它会通过执行对应的信号处理程序（signal handler）来进行响应。在shell脚本中同样可以发送、接收信号，并对其进行处理。KILL是用于终止进程的信号。像Ctrl+C、Ctrl+Z这种作业都属于信号。kill命令可用来向进程发送信号，而trap命令用来处理所接收的信号。

9.3.2 实战演练

列出所有可用的信号：

```
$ kill -l
```

该命令会打印出信号数（signal number）和信号名称。

终止一个进程：

```
$ kill PROCESS_ID_LIST
```

kill命令默认发出一个TERM信号。进程ID列表使用空格作为进程ID之间的定界符。

要通过kill命令向进程发送指定的信号，可以使用：

```
$ kill -s SIGNAL PID
```

参数SIGNAL要么是信号名称，要么是信号数。尽管可以指定很多信号用于不同的目的，我们经常用到的其实只有少数几个，具体如下所示。

- SIGHUP 1——对控制进程或终端进行挂起检测（hangup detection）。
- SIGINT 2——当按下Ctrl+C时发送该信号。
- SIGKILL 9——用于强行杀死进程。
- SIGTERM 15——默认用于终止进程。
- SIGTSTP 20——当按下Ctrl+Z时发送该信号。

我们经常要强行杀死进程，那么可以使用：

```
$ kill -s SIGKILL PROCESS_ID
```

或者

```
$ kill -9 PROCESS_ID
```

9.3.3 补充内容

让我们看看用于终止以及向进程发送信号的其他命令。

1. 杀死一组命令

kill命令以进程ID作为参数。在kill命令系列中还有其他命令可以接受命令名作为参数，并向对应的进程发送信号。

killall命令通过命令名终止进程：

```
$ killall process_name
```

通过名称向进程发送信号：

```
$ killall -s SIGNAL process_name
```

通过名称强行杀死进程：

```
$ killall -9 process_name
```

例如：

```
$ killall -9 gedit
```

通过名称以及所属用户名指定进程：

```
$ killall -u USERNAME process_name
```

如果需要在杀死进程前进行确认，可以使用killall的-i选项。

pkill命令和kill命令类似，不过默认情况下pkill接受的是进程名，而非进程ID。例如：

```
$ pkill process_name
```

```
$ pkill -s SIGNAL process_name
```

SIGNAL是信号数字。pkill不支持信号名称。

pkill提供了很多和kill相同的选项。要了解更多信息，请参阅pkill的命令手册。

2. 捕捉并响应信号

trap命令在脚本中用来为信号分配信号处理程序。一旦使用trap将某个函数分配给一个信号，那么当脚本运行时收到这个信号，对应于信号的函数就会开始执行。

命令语法如下：

```
trap 'signal_handler_function_name' SIGNAL LIST
```

SIGNAL LIST以空格分隔，它可以是信号数字或者信号名称。

写一个能够响应信号SIGINT的shell脚本：

```
#!/bin/bash
#文件名: sighandle.sh
#用途: 信号处理程序
function handler()
{
    echo Hey, received signal : SIGINT
}
echo My process ID is $$
# $$是一个特殊变量，它可以返回当前进程的进程ID
trap 'handler' SIGINT
#handler是信号SIGINT的信号处理程序的名称
while true;
```



```
do
    sleep 1
done
```

在终端运行这个脚本。当脚本运行时，如果按Ctrl+C，就会显示一条消息，这是通过执行与信号关联的信号处理程序实现的。Ctrl+C就是一个SIGINT信号。

通过使用一个无限循环while来保持进程一直运行。这样就可以使它能够响应另一个进程以异步方式发送的信号。用来保持进程一直处于活动状态的循环通常称为事件循环（event loop）。如果不使用无限循环，脚本在执行完所有语句之后就会终止。对于信号处理程序脚本而言，它必须等待并响应信号。

我们可以用kill命令以及脚本的进程ID向脚本发送信号：

```
$ kill -s SIGINT PROCESS_ID
```

上面的脚本在执行的时候会打印出 PROCESS_ID。或者，你也可以用ps命令找出脚本的进程ID。

如果没有为信号指定信号处理程序，那么将会调用操作系统默认分配的信号处理程序。一般来说，按下Ctrl+C会终止程序，这是因为操作系统提供的处理程序的默认行为。但是，在这里我们自行定义的信号处理程序指定了在接收到信号后所执行的特定行为。

通过trap命令，我们能够任意可用的信号（kill -l）定义处理程序，也可以为多个信号指定单个信号处理程序。

9.4 which、whereis、file、whatis 与平均负载

这则攻略旨在讲解我们通常会碰到的几个命令。掌握这些命令对于用户来说很有益处。

实战演练

让我们看看以下这些命令及其用例。

□ which

which命令用来找出某个命令的位置。我们在终端输入命令的时候无需知道对应的可执行文件位于何处。终端会在一组位置中查找这个命令，如果能够找到，那么就执行该可执行文件。这一组位置由环境变量PATH指定。例如：

```
$ echo $PATH
/usr/local/sbin:/usr/local/bin:/usr/sbin:/usr/bin:/sbin:/bin:/usr/games
```

我们可以导出并添加我们自己的命令搜索位置。例如，要将/home/slynux/bin添加到PATH中，可以使用以下命令：

```
$ export PATH=$PATH:/home/slynux/bin
# /home/slynux/bin被添加到PATH中
```

which命令输出作为参数的命令的所在位置。例如：

```
$ which ls
/bin/ls
```

❑ whereis

whereis与which命令类似，但是它不仅返回命令的路径，还能够打印出其对应的命令手册的位置以及命令源代码的路径（如果有的话）。例如：

```
$ whereis ls
ls: /bin/ls /usr/share/man/man1/ls.1.gz
```

❑ file

file命令是一个既有趣又使用频繁的命令。它用来确定文件的类型：

```
$ file FILENAME
```

该命令会打印出与该文件类型相关的细节信息。

例如：

```
$ file /bin/ls
/bin/ls: ELF 32-bit LSB executable, Intel 80386, version 1 (SYSV),
dynamically linked (uses shared libs), for GNU/Linux 2.6.15, stripped
```

❑ whatis

whatis命令会输出作为参数的命令的简短描述信息。这些信息是从命令手册中解析得来的。例如：

```
$ whatis ls
ls (1)                - list directory contents
```



apropos

有时候我们需要搜索和某个单词相关的命令是否存在。那么可以搜索包含该字符串命令的手册页。为此，我们使用：

```
apropos COMMAND
```

❑ 平均负载

平均负载（load average）是系统运行总负载量的一个重要参数。它指明了系统中可运行进程总量的平均值。平均负载由三个值来指定，第一个值指明了1分钟内的平均值，第二个值指明了5分钟内的平均值，第三个值指明了15分钟内的平均值。

这三个值可以通过运行uptime获得。例如：

```
$ uptime
12:40:53 up 6:16, 2 users, load average: 0.00, 0.00, 0.00
```

9.5 向用户终端发送消息

系统管理员可能需要向网络中所有主机上的所有用户或特定用户的终端发送消息。这则攻略将指导你如何完成这项任务。

9.5.1 新手上路

wall命令用来向所有当前登录用户的终端写入消息。它可以将消息传递给一台服务器中所有的当前登录用户或是多台分散主机中的用户。给所有的用户发送消息未必有用。不过我们可能需要给特定用户或终端发送消息。在Linux系统中，终端是作为设备存在的。因此那些打开的终端在/dev/pts/中都会有对应的设备节点文件。向特定的设备写入数据将会在对应的终端中显示出消息。

9.5.2 实战演练

向终端中所有的当前登录用户发送广播消息：

```
$ cat message | wall
```

或者

```
$ wall < message
Broadcast Message from slynux@slynux-laptop
(/dev/pts/1) at 12:54 ...
```

```
This is a message
```

消息概要(message outline)会显示谁(哪个用户、哪台主机)发送了这则消息。如果其他用户发送了消息，只有在“写入消息”选项启用的情况下该消息才会显示在当前终端中。在绝大多数发行版中，“写入消息”选项都是默认启用的。如果消息的发送者是超级用户，那么不管“写入消息”选项是否启用，消息都会显示出来。

要允许写入消息，可以使用：

```
$ mesg y
```

要禁止写入消息，可以使用：

```
$ mesg n
```

让我们写一个给指定用户终端发送消息的脚本：

```
#!/bin/bash
#文件名: message_user.sh
#用途: 用于向指定用户记录的终端发送信息的脚本
USER=$1
devices=`ls /dev/pts/* -l | awk '{ print $3,$9 }' | grep $USER | awk
'{ print $2 }'`
for dev in $devices;
do
    cat /dev/stdin > $dev
done
```

运行脚本：

```
./message_user.sh USERNAME < message.txt
# 通过stdin传递消息，USERNAME作为参数
```

输出如下:

```
$ cat message.txt
A message to slynux. Happy Hacking!
# ./message_user.sh slynux < message.txt
# 因为消息要发送给指定的用户, 因此要以超级用户身份运行message_user.sh
```

这时slynux的终端将接收到消息。

9.5.3 工作原理

目录 /dev/pts 中包含着对应于终端中每位当前登录用户的字符设备。我们可以通过查看设备文件的属主来得知谁登入了哪个终端。ls-l的输出包含了属主名以及设备路径。这些信息可以用awk提取。然后用grep单独提取对应于指定用户的行。用户名作为脚本的首个参数被存储在变量USER中, 然后生成给定用户的终端列表。for循环用来迭代每一个设备路径。/dev/stdin包含传递给当前进程的标准输入数据。通过读取 /dev/stdin, 就可以获取数据并将其重定向到对应的终端设备 (TTY)。因此消息就得以显示。

9.6 收集系统信息

从命令行中收集当前系统信息对于记录系统数据来说非常重要。各种系统信息包括主机名、内核版本、Linux发行版名称、CPU信息、内存信息、磁盘分区信息等。这则攻略将为你演示Linux中收集系统信息的不同方法。

实战演练

打印当前系统的主机名:

```
$ hostname
```

或者

```
$ uname -n
```

打印Linux内核版本、硬件架构等详细信息:

```
$ uname -a
```

打印内核发行版本:

```
$ uname -r
```

打印主机类型:

```
$ uname -m
```

打印出CPU的相关信息:

```
$ cat /proc/cpuinfo
```



获取处理器名称:

```
$ cat /proc/cpuinfo | head -n 5 | tail -1
```

cpuinfo的第5行包含处理器的名称。因此首先提取出前5行,然后再提取最后一行来打印处理器名称。

打印内存的详细信息:

```
$ cat /proc/meminfo
```

打印系统可用内存总量:

```
$ cat /proc/meminfo | head -1
MemTotal:      1026096 kB
```

列出系统的分区信息:

```
$ cat /proc/partitions
```

或者

```
$ fdisk -l
```

获取系统的详细信息:

```
$ lshw
```

9.7 用/proc 收集信息

在GNU/Linux操作系统中,/proc是一个位于内存中的伪文件系统(in-memory pseudo filesystem)。它的引入是为了提供一个可以从用户空间(user space)读取系统参数的接口。我们能够从中收集到大量的系统信息。来看proc文件系统的一些特性。

实战演练

如果查看/proc的话,你会发现很多文件和目录。其中的一些我们在本章的其他攻略中已经讲解过了。你可以对/proc中的文件和子目录进行cat来获取信息。所有内容都是易读的格式化文本。

系统中每一个运行的进程在/proc中都有一个对应的目录。进程的目录名和进程ID相同。

以Bash为例,它的进程ID是4295 (pgrep bash),那么将存在对应的目录 /proc/4295中。进程对应的目录包含了大量有关进程的信息。目录 /proc/PID中一些重要的文件如下所示。

□ environ——包含与进程相关联的环境变量。

使用cat /proc/4295/environ,可以显示所有传递给该进程的环境变量。

□ cwd——是一个到进程工作目录的符号链接。

```
$ readlink /proc/4295/exe
/bin/bash
```

□ fd——包含了由进程所使用的文件描述符。

9.8 用 cron 进行调度

通常我们需要安排脚本在某个时间或每隔一段时间来运行。GNU/Linux系统包含了各种用于调度任务的工具。cron就是其中之一，它通过守护进程cron使得任务能够以固定的时间间隔在系统后台自动运行。cron利用的是一个被称为“cron表”(cron table)的文件，这个文件中存储了需要执行的脚本或命令的调度列表以及执行时间。这个工具非常有用，一个常见的用例是设置在免费时段(free hour)(这是一些ISP提供免费服务，通常将时间设定在大多数人都已入睡的午夜)从Internet上进行下载。用户完全不需要在夜里熬红双眼等待下载。只需要编写一个cron条目，然后调度下载即可。你也可以安排当免费时段结束后自动断开Internet连接并关机。

9.8.1 新手上路

所有的GNU/Linux发行版都默认包含了cron调度工具。只要我们在cron表中写入条目，对应的命令就会在指定的时间执行。命令crontab用来添加调度条目。一次cron调度其实就是一段简单的文本。每一位用户都有自己的cron调度，通常这也被称为一个cron作业(cron job)。

9.8.2 实战演练

要想进行任务调度，我们得知道cron表的格式。一个cron作业指定了需要执行的脚本或命令的路径以及执行时间。cron表的每一个条目都由6部分组成，并按照下列顺序排列：

- 分钟 (0~59)
- 小时 (0~23)
- 天 (1~31)
- 月份 (1~12)
- 工作日 (0~6)
- 命令 (在指定时间执行的脚本或命令)

前5部分指定了某个命令实例所要执行的时间。当然还有其他选项也可用以指定调度时间。

星号(*)指定命令应该在每一个时间阶段执行。也就是说，如果*是写在cron作业中的小时字段中，那么命令就会每小时执行一次。与此类似，如果你希望在某个特定时段执行命令，那么就在对应的时间字段中指定时段，并用逗号分隔(例如要在第5分钟和第10分钟运行命令，那就在分钟字段中输入*5, 10*)。还有另一个不错的选项可以让我们以特定的时间间隔运行命令。在分钟字段使用*/5，可以每5分钟运行一次命令。这个技巧可以用在任何时间字段。一个cron表条目是由一行或多行cron作业组成的。cron表条目中的每一行都是一项作业。例如：

- 编写一个crontab条目的样例：

```
02 * * * * /home/slynux/test.sh
```

这个cron作业会在每天各小时的第2分钟执行脚本test.sh。

- 要在每天的第5、6、7小时执行脚本：

```
00 5,6,7 * * /home/slynux/test.sh
```

□ 在周日的每个小时执行脚本script.sh:

```
00 */12 * * 0 /home/slynux/script.sh
```

□ 在每天凌晨2点钟关闭计算机:

```
00 02 * * * /sbin/shutdown -h
```

现在,让我们看看如何调度一项cron作业。执行crontab命令进行调度的方法有很多种。

如果你手动运行crontab,用选项-e输入cron作业:

```
$ crontab -e
02 02 * * * /home/slynux/script.sh
```

当输入crontab -e后,会打开默认的文本编辑器(通常是vi)供用户输入cron作业并保存。这项cron作业将会在指定的时间被调度执行。

如果我们在脚本中调用crontab进行任务调度,那么有另外两种方法可供使用:

(1) 创建一个文本文件(例如task.cron),并写入cron作业。

然后将文件名作为命令参数,运行crontab:

```
$ crontab task.cron
```

(2) 通过下面的方法,我们可以在行内(inline)指定cron作业,而无需创建单独的文件。

例如:

```
crontab<<EOF
02 * * * * /home/slynux/script.sh
EOF
```

cron作业需要写在crontab<<EOF和EOF之间。

执行cron作业所使用的权限同执行crontab命令所使用的权限相同。如果你需要执行要求更高权限的命令,例如关闭计算机,那么就要以超级用户身份执行crontab。

在cron作业中指定的命令需要使用完整路径。这是因为执行cron作业时的环境与终端所使用的环境不同,因此环境变量PATH可能都没有设置。如果命令运行时需要设置某些环境变量,你应该明确设定出来。

9.8.3 补充内容

crontab命令还有其他选项。让我们看看其中一部分。

1. 指定环境变量

很多命令需要正确的设置环境变量才能够运行。我们可以在用户的cron表中插入一行变量赋值语句来设置环境变量。

例如,如果你使用的是代理服务器连接Internet,要调度某个需要使用Internet的命令,你就得设置HTTP代理环境变量http_proxy,可以用下面的方法来完成:

```
crontab<<EOF
http_proxy=http://192.168.03:3128
00 * * * * /home/slynux/download.sh
EOF
```

2. 查看cron表

我们可以用选项-l列出现有的cron表中的内容：

```
$ crontab -l
02 05 * * * /home/user/disklog.sh
```

crontab -l会列出当前用户cron表中的已有条目。

我们也可以通过选项-u指定用户名来查看其他用户的cron表：

```
$ crontab -l -u slynux
09 10 * * * /home/slynux/test.sh
```

当使用选项-u时，你应该使用超级用户以获取更高的权限。

3. 移除cron表

可以用选项-r移除当前用户的cron表：

```
$ crontab -r
```

要移除其他用户的cron表，可以使用：

```
# crontab -u slynux -r
```

这需要以超级用户身份获得更高的权限。

9.9 从 Bash 中读写 MySQL 数据库

MySQL是一款应用广泛的数据库系统。MySQL数据库通常是被以PHP、Python、C++等语言编写的应用程序用作存储系统。从shell脚本中访问并操作MySQL数据库很有意思。我们可以编写脚本将文本文件或CSV（Comma Separated Value）的内容写入数据表，与MySQL进行交互来读取并处理数据。例如，我们可以通过从shell脚本中执行查询语句来读取存储在留言板（guestbook）程序的数据库中所有电子邮件地址。在这则攻略中，我们会看到如何从Bash中读写MySQL数据库。为了便于演示，这里给出一个示例问题：

我有一个包含多个系的学生详细信息的CSV文件。我需要将文件的内容插入到一个数据表中。要保证为每一个系生成一个单独的排名列表。

9.9.1 新手上路

要处理MySQL数据库，系统中必须安装mysql-server和mysql-client软件包。Linux发行版并没有默认包含这些工具。由于MySQL要使用用户名和密码进行认证，因此你得有用户名和密码才能运行脚本。

9.9.2 实战演练

前面提出的问题可以用sort、awk等Bash工具解决，或者用一个SQL数据库的数据表也可以搞定。我们接下来要编写3个脚本，分别用于创建数据库及数据表、向数据表中插入学生数据、从数据表中读取并显示处理过的数据。

创建数据库及数据表的脚本如下：

```
#!/bin/bash
#文件名: create_db.sh
#用途: 创建MySQL数据库和数据表

USER="user"
PASS="user"

mysql -u $USER -p$PASS <<EOF 2> /dev/null
CREATE DATABASE students;
EOF

[ $? -eq 0 ] && echo Created DB || echo DB already exist
mysql -u $USER -p$PASS students <<EOF 2> /dev/null
CREATE TABLE students(
id int,
name varchar(100),
mark int,
dept varchar(4)
);
EOF

[ $? -eq 0 ] && echo Created table students || echo Table students
already exist

mysql -u $USER -p$PASS students <<EOF
DELETE FROM students;
EOF
```

将数据插入数据表的脚本如下：

```
#!/bin/bash
#文件名: write_to_db.sh
#用途: 从CSV中读取数据并写入MySQLdb

USER="user"
PASS="user"

if [ $# -ne 1 ];
then
echo $0 DATAFILE
echo
exit 2
fi
data=$1

while read line;
do
oldIFS=$IFS
IFS=,
values=($line)
values[1]="\"echo ${values[1]} | tr ' ' '#' \"\"
values[3]="\"echo ${values[3]}\"\"

query='echo ${values[@]} | tr ' #' , , '
IFS=$oldIFS
```

```
mysql -u $USER -p$PASS students <<EOF
INSERT INTO students VALUES($query);
EOF
```

```
done< $data
echo Wrote data into DB
```

查询数据库的脚本如下:

```
#!/bin/bash
#文件名: read_db.sh
#用途:从数据库中读取数据

USER="user"
PASS="user"

depts=`mysql -u $USER -p$PASS students <<EOF | tail -n +2
SELECT DISTINCT dept FROM students;
EOF`

for d in $depts;
do
    echo Department : $d

    result=`mysql -u $USER -p$PASS students <<EOF
SET @i:=0;
SELECT @i:=@i+1 as rank,name,mark FROM students WHERE dept="$d" ORDER
BY mark DESC;
EOF`

    echo "$result"
    echo
done
```

作为输入的CSV文件 (studentdata.csv) 中的数据如下:

```
1,Navin M,98,CS
2,Kavya N,70,CS
3,Nawaz O,80,CS
4,Hari S,80,EC
5,Alex M,50,EC
6,Neenu J,70,EC
7,Bob A,30,EC
8,Anu M,90,AE
9,Sruthi,89,AE
10,Andrew,89,AE
```

按照以下顺序执行脚本:

```
$ ./create_db.sh
Created DB
Created table students

$ ./write_to_db.sh studentdat.csv
Wrote data into DB

$ ./read_db.sh
```



```

Department : CS
rank name mark
1 Navin M 98
2 Nawas O 80
3 Kavya N 70

```

```

Department : EC
rank name mark
1 Hari S 80
2 Neenu J 70
3 Alex M 50
4 Bob A 30

```

```

Department : AS
rank name mark
1 Anu M 90
2 Sruthi 89
3 Andrew 89

```

9.9.3 工作原理

我们现在来逐个讲解上面的脚本。第一个脚本`create_db.sh`用来创建数据库`students`，并在其中创建数据表`students`。我们需要MySQL的用户名和密码来访问或修改数据库中的内容。变量`USER`和`IPASS`用来存储用户名和密码。`mysql`命令用于对MySQL进行操作。`mysql`命令可以用`-u`指定用户名，用`-pPASSWORD`指定密码，其他命令参数是数据库名。如果将数据库名作为`mysql`命令的参数，那么就将使用该数据库，否则我们必须用`use database_name`明确地指定SQL查询语句使用哪一个数据库进行查询。`mysql`命令通过标准输入（`stdin`）接受查询。通过`stdin`提供多行输入的简便方法是使用`<<EOF`。出现在`<<EOF`和`EOF`之间的文本被作为`mysql`的标准输入。在`CREATE DATABASE`语句中，为了避免显示错误信息，我们将`stderr`重定向到`/dev/null`。同样，在创建数据表时，我们也将`stderr`重定向到`/dev/null`，以忽略可能出现的任何错误。然后我们用退出状态变量`$?`来检查`mysql`命令的退出状态，以获知是否已经存在同名的数据库或数据表。如果已经存在，则会显示出一条提示信息，否则，就进行创建。

接下来的脚本`write_to_db.sh`接受包含学生数据的CSV文件名。我们用`while`循环读取CSV文件的每一行，所以在每次迭代中都会接收到一行以逗号分隔的数值。然后我们需要将行内的数值放入SQL查询语句中。要实现这个目的，最简单的方法是用数组存储CSV文件行中的数据项。我们知道数组赋值的形式为`array=(val1 val2 val3)`，其中内部字段分隔符（IFS）是空格。我们的文本行用逗号分隔数值，因此只需要将IFS修改成逗号（IFS=,），我们就可以轻松地赋值给数组。文本行中以逗号分隔的数据项分别是`id`、`name`、`mark`和`department`。`id`和`mark`是整数，而`name`和`department`是字符串（字符串必须进行引用）。`name`中也可以包含空格。这样一来就和IFS产生了冲突。因此我们应该将`name`中的空格替换成其他字符（#），在构建查询语句时再替换回来。为了引用字符串，数组中的值要加上`\`作为前缀和后缀。`tr`用来将`name`中的空格替换成#。最后通过将空格替换成逗号，将#替换成空格来构造出查询语句并进行查询。

第三个脚本`read_db.sh`用来查找各系并打印出每个系的学生排名列表。第一个查询用来找出各系的名称。我们用一个`while`循环对每个系进行迭代，然后进行查询并按照成绩从高到低的顺

序显示学生的详细信息。SET @i:=0是一个SQL构件 (SQL construct)，用来设置变量i=0。在每一行中，变量i都会进行增加并作为学生排名来显示。

9.10 用户管理脚本

GNU/Linux是一个多用户操作系统。多个用户可以同时登录并执行多种操作。而有一些管理任务会涉及用户管理，这些任务包括为用户设置默认shell、禁用某个账户、禁用某个shell账户、添加新用户、删除用户、设置密码、设置账户有效期等。这则攻略旨在编写一个可以处理这类任务的用户管理工具。

9.10.1 实战演练

让我们看看这个用户管理脚本：

```
#!/bin/bash
#文件名: user_admin.sh
#用途: 用户管理工具

function usage()
{
    echo Usage:
    echo Add a new user
    echo $0 -adduser username password
    echo
    echo Remove an existing user
    echo $0 -deluser username
    echo
    echo Set the default shell for the user
    echo $0 -shell username SHELL_PATH
    echo
    echo Suspend a user account
    echo $0 -disable username
    echo
    echo Enable a suspended user account
    echo $0 -enable username
    echo
    echo Set expiry date for user account
    echo $0 -expiry DATE
    echo
    echo Change password for user account
    echo $0 -passwd username
    echo
    echo Create a new user group
    echo $0 -newgroup groupname
    echo
    echo Remove an existing user group
    echo $0 -delgroup groupname
    echo
    echo Add a user to a group
    echo $0 -addgroup username groupname
```



```

echo
echo Show details about a user
echo $0 -details username
echo
echo Show usage
echo $0 -usage
echo

exit

}

if [ $UID -ne 0 ];
then
    echo Run $0 as root.
    exit 2
fi

case $1 in

-adduser) [ $# -ne 3 ] && usage ; useradd $2 -p $3 -m ;;
-deluser) [ $# -ne 2 ] && usage ; deluser $2 --remove-all-files;;
-shell)    [ $# -ne 3 ] && usage ; chsh $2 -s $3 ;;
-disable)  [ $# -ne 2 ] && usage ; usermod -L $2 ;;
-enable)   [ $# -ne 2 ] && usage ; usermod -U $2 ;;
-expiry)   [ $# -ne 3 ] && usage ; chage $2 -E $3 ;;
-passwd)   [ $# -ne 2 ] && usage ; passwd $2 ;;
-newgroup) [ $# -ne 2 ] && usage ; addgroup $2 ;;
-delgroup) [ $# -ne 2 ] && usage ; delgroup $2 ;;
-addgroup) [ $# -ne 3 ] && usage ; addgroup $2 $3 ;;
-details)  [ $# -ne 2 ] && usage ; finger $2 ; chage -l $2 ;;
-usage)    usage ;;
*)         usage ;;

esac

```

输出样例如下:

```

# ./user_admin.sh -details test
Login: test                               Name:
Directory: /home/test                     Shell: /bin/sh
Last login Tue Dec 21 00:07 (IST) on pts/1 from localhost
No mail.
No Plan.
Last password change                      : Dec 20, 2010
Password expires                          : never
Password inactive                        : never
Account expires                          : Oct 10, 2010
Minimum number of days between password change : 0
Maximum number of days between password change : 99999
Number of days of warning before password expires : 7

```

9.10.2 工作原理

脚本user_admin.sh可以用来执行多项用户管理任务。你可以参考usage()中的内容来学习这个脚本的用法。当用户给出的参数不正确或使用-usage选项时,函数usage()用来显示脚本不同

选项的使用方法。case语句用来匹配命令参数，并根据参数执行对应的命令。脚本user_admin.sh合法的命令选项是：-adduser, -deluser, -shell, -disable, -enable, expiry, -passwd, -newgroup, -delgroup, -addgroup, -details和-usage。如果匹配了*)分支，那就意味着用户输入了错误的选项，因此要调用usage()。对于每一个匹配分支，我们使用[\$# -ne 3] && usage。它用来检测参数个数。如果命令参数个数不等于要求的数量，则调用函数usage()并退出脚本。要运行用户管理命令，需要以超级用户身份执行脚本。因此要检查用户ID是否为0（超级用户的用户ID是0）。如果用户ID非0，则表明脚本不是以超级用户身份执行的。因此显示出要求以超级用户身份运行脚本的提示信息并退出。

下面来逐个讲解每个选项。

❑ -useradd

useradd命令可以用来创建新用户。命令语法如下：

```
useradd USER -p PASSWORD
```

选项-m用来创建home目录。

也可以用选项-c FULLNAME提供用户的全名。

❑ -deluser

deluser命令用来删除用户。命令语法如下：

```
deluser USER
```

--remove-all-files用来删除与用户相关的所有文件，包括home目录。

❑ -shell

chsh命令用来修改用户的默认shell。命令语法如下：

```
chsh USER -s SHELL
```

❑ -disable和-enable

usermod命令用来处理和用户账户相关的若干属性信息。

usermod -L USER和usermod -U USER分别用来锁定和解锁用户账户。

❑ -expiry

chage命令用来处理用户账户的过期信息。命令语法如下：

```
chage -E DATE
```

其他选项如下：

- -m MIN_DAYS (将更改密码的最小天数修改成MIN_DAYS)；
- -M MAX_DAYS (设置密码有效的最大天数)；
- -W -WARN_DAYS (设置在前几天提醒需要更改密码)。

❑ -passwd

passwd命令用来更改用户密码。命令语法如下：

```
passwd USER
```

命令会提示输入新的密码。

❑ `-newgroup`和`addgroup`

`addgroup`命令会为系统添加一个新的用户组。命令语法如下：

```
addgroup GROUP
```

要将已有的用户添加到一个组，可以使用：

```
addgroup USER GROUP
```

❑ `-delgroup`

`delgroup`命令会删除一个用户组。命令语法如下：

```
delgroup GROUP
```

❑ `-details`

`finger USER`命令会显示用户信息，这包括用户的home目录、上一次登录的时间、默认shell等。`chage -l`命令会显示用户账户的过期信息。

9.11 图像文件的批量缩放及格式转换

我们大家都会使用数码相机，也会从Internet上获取数码相机照片。如果需要处理大量图像文件，我们可以轻松地用脚本批量处理。这通常会涉及的任务是调整照片的大小，有时也需要转换图像格式（例如，将JPEG格式转换成PNG格式）。当我们从数码相机中获取到照片时，大分辨率的图片通常体积都比较大，我们可能需要减少图片的大小，以便于存储以及通过电子邮件发送给他人。因此我们就需要对其进行调整来降低分辨率。这则攻略将讨论如何用脚本管理图像。

9.11.1 新手上路

Imagemagick是一款出色的图像处理工具，它包含丰富的选项，能够处理多种图像格式。大多数GNU/Linux发行版并没有安装Imagemagick。你得自己手动安装这个软件包。`convert`是我们经常要使用的命令。

9.11.2 实战演练

将一种图像格式转换为另一种图像格式：

```
$ convert INPUT_FILE OUTPUT_FILE
```

例如：

```
$ convert file1.png file2.png
```

我们可以通过指定缩放比或输出图像的宽度和高度来将图像调整到所需要的大小。指定WIDTH（宽度）或HEIGHT（高度）来缩放图像：

```
$ convert image.png -resize WIDTHxHEIGHT image.png
```

例如:

```
$ convert image.png -resize 1024x768 image.png
```

必须提供WIDTH或HEIGHT, 这样才能够使脚本自动计算其他数值, 以便于在保留图像比例的同时进行缩放。

```
$ convert image.png -resize WIDTHx image.png
```

例如:

```
$ convert image.png -resize 1024x image.png
```

指定百分比缩放图像:

```
$ convert image.png -resize "50%" image.png
```

让我们看一个用于图像管理的脚本:

```
#!/bin/bash
# 文件名: image_help.sh
# 用途: 图像管理脚本

if [ $# -ne 4 -a $# -ne 6 -a $# -ne 8 ];
then
    echo Incorrect number of arguments
    exit 2
fi

while [ $# -ne 0 ];
do
    case $1 in
        -source) shift; source_dir=$1 ; shift ;;
        -scale) shift; scale=$1 ; shift ;;
        -percent) shift; percent=$1 ; shift ;;
        -dest) shift ; dest_dir=$1 ; shift ;;
        -ext) shift ; ext=$1 ; shift ;;
        *) echo Wrong parameters; exit 2 ;;
    esac;
done

for img in `echo $source_dir/*` ;
do
    source_file=$img
    if [[ -n $ext ]];
    then
        dest_file=${img%.*}.$ext
    else
        dest_file=$img
    fi

    if [[ -n $dest_dir ]];
    then
        dest_file=${dest_file##*/}
        dest_file="$dest_dir/$dest_file"
    fi

    if [[ -n $scale ]];
```



```

then
    PARAM="-resize $scale"
elif [[ -n $percent ]];
then
    PARAM="-resize $percent%"
fi

echo Processing file : $source_file
convert $source_file $PARAM $dest_file

done

```

下面是输出样例，将目录sample_dir中的图像调整到原来的20%：

```

$ ./image_help.sh -source sample_dir -percent 20%
Processing file :sample/IMG_4455.JPG
Processing file :sample/IMG_4456.JPG
Processing file :sample/IMG_4457.JPG
Processing file :sample/IMG_4458.JPG

```

将图像宽度调整到1024：

```

$ ./image_help.sh -source sample_dir -scale 1024x

```

把-ext png加入上面的命令，使文件格式转换成PNG。

将文件缩放或转换到指定的目录：

```

$ ./image_help.sh -source sample -scale 50% -ext png -dest newdir
# newdir作为目的目录

```

9.11.3 工作原理

上面的脚本image_help.sh可以接受多个命令行参数，例如-source、-percent、-scale、-ext和-dest等。每个选项的简短描述如下。

- -source用于指定图像源目录。
- -percent用于指定缩放比例，-scale用于指定缩放宽度与高度。
- -percent与-scale不能同时出现，只能使用其中之一。
- -ext用于指定目标文件格式。-ext是一个可选的选项。如果没有指定，那么不执行格式转换。
- -dest为缩放或转换格式后的文件指定目的目录。该选项也是可选的。如果没有指定，目的目录则和源目录相同。脚本的第一步就是检查命令行参数的数量是否正确，可以出现的参数数量分别是4、6或8。

借助while循环和case语句，我们将命令行参数解析到对应的变量。\$#是一个特殊的变量，它可以返回参数的数量。shift命令每执行一次，就将命令行参数向左移动一个位置，这样我们就需要使用变量\$1、\$2、\$3等，而只用一个\$1来就可以对命令参数逐个进行访问了。case语句用来匹配\$1的值，就像C语言中的switch语句一样。如果匹配了某个case分支，就执行对应的语句。每一个case分支都以;;作为结尾。一旦将所有的参数都解析到变量percent、scale、source_dir、ext和dest_dir中，就是用for循环对源目录中的每一个文件进行迭代，并执行对应的转换操作。

如果变量`ext`已定义（也就是说`-ext`作为命令参数出现），就将目标文件的扩展名从`source_file.extension`更改为`source_file.$ext`。接下来检查是否提供了`-dest`选项。如果指定了目的目录，则使用文件名切片将源路径中的目录替换成目的目录，从而形成目的文件路径。然后构造出`convert`命令的参数，用以执行缩放（`-resize widthx` 或 `-resize perc%`）。参数构造完毕之后，用对应的参数执行`convert`命令。

9.11.4 参考

2.11节讲解了如何提取部分文件名。

