The bash shell and UNIX commands

Chris Paciorek

2024-09-03

Table of contents

| 1. | Shell basics | 1 |
|------------|---|---|
| 2. | Using the bash shell | 2 |
| 3. | bash shell examples | 2 |
| 4. | bash shell challenges | 6 |
| | 4.1 First challenge | 6 |
| | 4.2 Second challenge | 6 |
| | 4.3 Third challenge | 7 |
| | 4.4 Fourth challenge | 7 |
| | 4.5 Fifth challenge | |
| | 4.6 Sixth challenge | 7 |
| 5 . | Regular expressions | 8 |
| | Versions of regular expressions | 8 |
| | General principles for working with regex | 9 |
| | Challenge problem | 9 |

PDF

Reference:

• Newham and Rosenblatt, Learning the bash Shell, 2nd ed.

1. Shell basics

The shell is the interface between you and the UNIX operating system.

I'll use 'UNIX' to refer to the family of operating systems that descend from the path-breaking UNIX operating system developed at AT&T's Bell Labs in the 1970s. These include MacOS and various flavors of Linux (e.g., Ubuntu, Debian, CentOS, Fedora).

When you are working in a terminal window (i.e., a window providing the command line interface), you're interacting with a shell. From the shell you can run UNIX commands such as cp, 1s, grep, etc.

(as well as start various applications).

Here's a graphical representation of how the shell relates to various programs, commands, and the operating system.

There are multiple shells (sh, bash, zsh, csh, tcsh, ksh). We'll assume usage of bash, as this is a very commonly-used shell in Linux, plus was the default for Mac OS X until Catalina (zsh, very similar to bash, is now the default), the SCF machines, and the UC Berkeley campus cluster (Savio). All of the various shells allow you to run UNIX commands.

For your work on this unit, either bash on a Linux machine, the older version of bash on MacOS, or zsh on MacOS (or Linux) are fine. I'll probably demo everything using bash on a Linux machine, and there are some annoying differences from the older bash on MacOS that may be occasionally confusing (in particular the options to various commands can differ on MacOS).

The Windows PowerShell and old cmd.exe/DOS command interpreter both provide a command-line interface on Windows, but not a UNIX-based one and not interfaces that will be considered here.

UNIX shell commands are designed to each do a specific task really well and really fast. They are modular and composable, so you can build up complicated operations by combining the commands. These tools were designed decades ago, so using the shell might seem old-fashioned, but the shell still lies at the heart of modern scientific computing. By using the shell you can automate your work and make it reproducible. And once you know how to use it, you'll find that enter commands quickly and without a lot of typing.

2. Using the bash shell

For this Unit, we'll rely on the bash shell tutorial for the details of how to use the shell. We won't cover the page on Managing Processes. For the moment, we won't cover the page on Regular Expressions, but when we talk about string processing and regular expressions in Unit 5, we'll come back to that material.

3. bash shell examples

Here we'll work through a few examples to start to give you a feel for using the bash shell to manage your workflows and process data.

First let's get the files from GitHub to have a set of files we can do interesting things with.

```
git clone https://github.com/berkeley-stat243/fall-2024
```

One important note is that most of the shell commands that work with data inside files (in contrast to commands like 1s and cd) work only with text files and not binary files. Also the commands operate on a line-by-line basis.

Our first mission is some basic manipulation of a data file. Suppose we want to get a sense for the number of weather stations in different states using the *coop.txt* file.

```
cd fall-2024/data
gzip -cd coop.txt.gz | less
```

```
gunzip coop.txt.gz
cut -b50-70 coop.txt | less
cut -b60-61 coop.txt | uniq
cut -b60-61 coop.txt | sort | uniq
cut -b60-61 coop.txt | sort | uniq -c
## Do it all in one line with no change to the original file:
gzip -cd coop.txt.gz | cut -b60-61 coop.txt | sort | uniq -c
```

I could have done that in R or Python, but it would have required starting the program up and reading all the data into memory.

If you feel that manually figuring out the position of the state field is inconsistent with our emphasis on programmatic workflows, see the fifth challenge below.

Our second mission: how can I count the number of fields in a CSV file programmatically?

```
tail -n 1 cpds.csv | grep -o ',' | wc -l
nfields=$(tail -n 1 cpds.csv | grep -o ',' | wc -l)

nfields=$((${nfields}+1))
echo $nfields

## Alternatively, we can use `bc`.
nfields=$(echo "${nfields}+1" | bc)
```

Trouble-shooting: How could the syntax above get the wrong answer?

Extension: We could write a function that can count the number of fields in any file.

Extension: How could I see if all of the lines have the same number of fields?

We'll work on Challenge #1 here.

Our third mission: was the requests package in the five most recently modified Quarto Markdown files in the units directory?

```
cd ../units
grep -l 'import requests' unit4-goodPractices.qmd
ls -tr *.qmd
## If unit4-goodPractices.qmd is not amongst the 5 most recently used,
## let's artificially change the timestamp so it is recently used.
touch unit4-goodPractices.qmd

ls -tr *.qmd | tail -n 5 | grep requests
ls -tr *.qmd | tail -n 5 | grep "unit4-goodPractices"
ls -tr *.qmd | tail -n 5 | xargs grep 'import requests'
ls -tr *.qmd | tail -n 5 | xargs grep -l 'import requests'
```

Notice that man tail indicates it can take input from a FILE or from stdin. Here it uses stdin, so

it is gives the last five lines of the output of ls, not the last five lines of the files indicated in that output.

man grep also indicates it can take input from a FILE or from stdin. However, we want grep to operate on the content of the files indicated in stdin. So we use xargs to convert stdin to be recognized as arguments, which then are the FILE inputs to grep.

An alternative to xargs is to embed the ls invocation, using \$(), to explicitly pass the file names as the argument to grep:

```
grep -l 'import requests' $(ls -tr *.R | tail -n 5)
```

Here are some of the ways we can pass information from a command to somewhere else:

- Piping allows us to pass information from one command to another command via stdout to stdin.
- \$() allows us to store the result of a command in a variable.
 - also used to create a temporary variable to pass the output from one command as an option or argument (e.g., the FILE argument) of another command
- File redirection operators such as > and >> allow us to pass output from a command into a file.

We'll work on Challenge #2 here.

Our fourth mission: write a function that will move the most recent n files in your Downloads directory to another directory.

In general, we want to start with a specific case, and then generalize to create the function.

Note that the quotes deal with cases where a file has a space in its name.

If we wanted to handle multiple files, we could do it with a loop:

Side note: if we were just moving files from the current working directory and with files without spaces in their names, it should be possible to use tail -n \${1} without the loop.

We'll work on Challenge #3 here.

We'll work on Challenge #4 here.

Our fifth mission: automate the process of determining what Python packages are used in all of the qmd code chunks here and install those packages on a new machine.

```
grep import *.qmd
grep --no-filename import *.qmd
grep --no-filename "^import" *.qmd
grep --no-filename "^import " *.qmd
grep --no-filename "^import " *.qmd | sort | uniq
grep --no-filename "^import " *.qmd | cut -d'#' -f1
grep --no-filename "^import " *.qmd | cut -d'#' -f1 | sed "s/as .*//"
grep --no-filename "^import " *.qmd | cut -d'#' -f1 | \
                   sed "s/as .*//" | sed "s/import //" > tmp.txt
sed s/\sqrt{n/g} tmp.txt | sed s/\sqrt{g} | sort | uniq | tee requirements.txt
## Note: on a Mac, use 's/,/\\n/g'
## See https://superuser.com/questions/307165/newlines-in-sed-on-mac-os-x
echo "There are $(wc -1 requirements.txt | cut -d' ' -f1) \
unique packages we will install."
## Note: on Linux, wc -l puts the number as the first characters of the output.
## On a Mac, there may be a bunch of spaces preceding the number, so try this:
## echo "There are $(wc -l libs.txt | tr -s ' ' | cut -d' ' -f2) \
## unique packages we will install."
pip install -r requirements.txt
# or use Mamba/Conda
```

You probably wouldn't want to use this code to accomplish this task in reality - you would want to see if there are packages that can accomplish this. (In the R ecosystem, the renv and packrat packages do this for R projects.) The main point was to illustrate how one can quickly hack together some code to do fairly complicated tasks.

Our sixth mission: suppose I've accidentally started a bunch of jobs (perhaps with a for loop in bash!) and need to kill them. (This example uses syntax from the Managing Processes page of the bash tutorial, so it goes beyond what you were asked to read for this Unit.)

```
# Use a 'here document':
cat > job.py << EOF
import time
time.sleep(1e5)
EOF</pre>
```

```
# alternatively:
echo -e "import time\ntime.sleep(1e5)" > job.py

nJobs=30
for (( i=1; i<=${nJobs}; i++ )); do
    python job.py > job-${i}.out &
done

# on Linux:
ps -o pid,pcpu,pmem,user,cmd -C python
ps -o pid,pcpu,pmem,user,cmd,start_time --sort=start_time -C python | tail -n 30
ps -o pid --sort=start_time -C python | tail -n ${nJobs} | xargs kill

# on a Mac:
ps -o pid,pcpu,pmem,user,command | grep python
# not clear how to sort by start time
ps -o pid,command | grep python | cut -d' ' -f1 | tail -n ${nJobs} | xargs kill
```

4. bash shell challenges

4.1 First challenge

Consider the file cpds.csv. How would you write a shell command that returns "There are 8 occurrences of the word 'Belgium' in this file.", where '8' should instead be the correct number of times the word occurs.

Extra: make your code into a function that can operate on any file indicated by the user and any word of interest.

4.2 Second challenge

Consider the data in the RTADataSub.csv file. This is a subset of data giving freeway travel times for segments of a freeway in an Australian city. The data are from a kaggle.com competition. We want to try to understand the kinds of data in each field of the file. The following would be particularly useful if the data were in many files or the data were many gigabytes in size.

- 1. First, take the fourth column. Figure out the unique values in that column.
- 2. Next, automate the process of determining if any of the values are non-numeric so that you don't have to scan through all of the unique values looking for non-numbers. You'll need to look for the following regular expression pattern [^0-9], which is interpreted as NOT any of the numbers 0 through 9.

Extra: do it for all the fields, except the first one. Have your code print out the result in a human-readable way understandable by someone who didn't write the code. For simplicity, you can assume you know the number of fields.

4.3 Third challenge

- 1. For Belgium, determine the minimum unemployment value (field #6) in cpds.csv in a programmatic wav.
- 2. Have what is printed out to the screen look like "Belgium 6.2".
- 3. Now store the unique values of the countries in a variable, first stripping out the quotation marks.
- 4. Figure out how to automate step 1 to do the calculation for all the countries and print to the screen.
- 5. How would you instead store the results in a new file?

4.4 Fourth challenge

Let's return to the RTADataSub.csv file and the issue of missing values.

- 1. Create a new file without any rows that have an 'x' (which indicate a missing value).
- 2. Turn the code into a function that also prints out the number of rows that are being removed and that sends its output to stdout so that it can be used with piping.
- 3. Now modify your function so that the user could provide the missing value string and the input filename.

4.5 Fifth challenge

Consider the coop.txt weather station file.

Figure out how to use grep to tell you the starting position of the state field. Hints: search for a known state-country combination and figure out what flags you can use with grep to print out the "byte offset" for the matched state.

Use that information to automate the first mission where we extracted the state field using cut. You'll need to do a bit of arithmetic using shell commands.

4.6 Sixth challenge

Here's an advanced one - you'll probably need to use sed, but the brief examples of text substitution in the using bash tutorial (or in the demos above) should be sufficient to solve the problem.

Consider a CSV file that has rows that look like this:

```
1, "America, United States of", 45, 96.1, "continental, coastal" 2, "France", 33, 807.1, "continental, coastal"
```

While Pandas would be able to handle this using read_csv(), using cut in UNIX won't work because of the commas embedded within the fields. The challenge is to convert this file to one that we can use cut on, as follows.

Figure out a way to make this into a new delimited file in which the delimiter is not a comma. At least one solution that will work for this particular two-line dataset does not require you to use regular expressions, just simple replacement of fixed patterns.

5. Regular expressions

Regular expressions ("regex") are a *domain-specific language* for finding and manipulating patterns of characters and are a key tool used in UNIX commands such as grep, sed, and awk as well as in scripting languages such as Python and R.

For the moment we'll focus on learning regular expression syntax in the context of the shell, but in Unit 5, we'll also use regular expressions within Python using the re package.

The basic idea of regular expressions is that they allow us to find matches of strings or patterns in strings, as well as do substitution. Regular expressions are good for tasks such as:

- extracting pieces of text;
- creating variables from information found in text;
- cleaning and transforming text into a uniform format; and
- mining text by treating documents as data.

Please see the bash shell tutorial for a description of regular expressions. I'll assign a small set of regex problems due as an assignment, and we'll talk through those problems in class to explore the regex syntax.

Other resources include:

- Here's a website where you can interactively test regular expressions on example strings.
- Duncan Temple Lang (UC Davis Statistics) has written a nice tutorial covering regular expressions, illustrated in R.
- Sections 9.9 and 11 of Paul Murrell's book
- The back/second page of RStudio's stringr cheatsheet has a cheatsheet on regular expressions.

In addition to the regular expression functionality we'll cover, there is a lot more advanced functionality we won't cover, such as callbacks, named groups, recursion, word boundaries, and lookaround assertions.

Versions of regular expressions

One thing that can cause headaches is differences in version of regular expression syntax used. As discussed in man grep, extended regular expressions are standard, with basic regular expressions providing less functionality and Perl regular expressions additional functionality.

The bash shell tutorial provides a full documentation of the extended regular expressions syntax, which we'll focus on here. This syntax should be sufficient for most usage and should be usable in Python and R, but if you notice something funny going on, it might be due to differences between the regular expressions versions.

- In bash, grep -E (or egrep) enables use of the extended regular expressions, while grep -P enables Perl-style regular expressions.
- In Python, the re package provides syntax "similar to" Perl.
- In R, stringr provides ICU regular expressions (see help(regex)), which are based on Perl regular expressions.

More details about Perl regular expressions can be found in the regex Wikipedia page.

General principles for working with regex

The syntax is very concise, so it's helpful to break down individual regular expressions into the component parts to understand them. As Murrell notes, since regex are their own language, it's a good idea to build up a regex in pieces as a way of avoiding errors just as we would with any computer code. re.findall in Python and str_detect in R's stringr, as well as regex101.com are particularly useful in seeing what was matched to help in understanding and learning regular expression syntax and debugging your regex. As with many kinds of coding, I find that debugging my regex is usually what takes most of my time.

Challenge problem

Challenge: Let's think about what regex syntax we would need to detect any number, integer- or real-valued. Let's start from a test-driven development perspective of writing out test cases including:

- various cases we want to detect,
- · various tricky cases that are not numbers and we don't want to detect, and
- "corner cases" tricky (perhaps unexpected) cases that might trip us up.