

< (<http://swcarpentry.github.io/shell-novice/05-loop/>)

The Unix Shell (<http://swcarpentry.github.io/shell-novice/>)

Shell Scripts

> (<http://swcarpentry.github.io/shell-novice/07-find/>)

? Overview

Teaching: 15 min

Exercises: 0 min

Questions

- How can I save and re-use commands?

Objectives

- Write a shell script that runs a command or series of commands for a fixed set of files.
- Run a shell script from the command line.
- Write a shell script that operates on a set of files defined by the user on the command line.
- Create pipelines that include shell scripts you, and others, have written.

We are finally ready to see what makes the shell such a powerful programming environment. We are going to take the commands we repeat frequently and save them in files so that we can re-run all those operations again later by typing a single command. For historical reasons, a bunch of commands saved in a file is usually called a **shell script**, but make no mistake: these are actually small programs.

Let's start by going back to `molecules/` and putting the following line into a new file, `middle.sh`:

```
$ cd molecules
$ nano middle.sh
```

The command `nano middle.sh` opens the file `middle.sh` within the text editor “nano” (which runs within the shell). If the file does not exist, it will be created. We can use the text editor to directly edit the file—we'll simply insert the following line:

```
head -n 15 octane.pdb | tail -n 5
```

This is a variation on the pipe we constructed earlier: it selects lines 11-15 of the file `octane.pdb`. Remember, we are *not* running it as a command just yet: we are putting the commands in a file.

Then we save the file (using `Ctrl-O`), and exit the text editor (using `Ctrl-X`). Check that the directory `molecules` now contains a file called `middle.sh`.

Once we have saved the file, we can ask the shell to execute the commands it contains. Our shell is called `bash`, so we run the following command:

```
$ bash middle.sh
```

ATOM	9	H	1	-4.502	0.681	0.785	1.00	0.00
ATOM	10	H	1	-5.254	-0.243	-0.537	1.00	0.00
ATOM	11	H	1	-4.357	1.252	-0.895	1.00	0.00
ATOM	12	H	1	-3.009	-0.741	-1.467	1.00	0.00
ATOM	13	H	1	-3.172	-1.337	0.206	1.00	0.00

Sure enough, our script's output is exactly what we would get if we ran that pipeline directly.

✈ Text vs. Whatever

We usually call programs like Microsoft Word or LibreOffice Writer “text editors”, but we need to be a bit more careful when it comes to programming. By default, Microsoft Word uses .docx files to store not only text, but also formatting information about fonts, headings, and so on. This extra information isn't stored as characters, and doesn't mean anything to tools like head: they expect input files to contain nothing but the letters, digits, and punctuation on a standard computer keyboard. When editing programs, therefore, you must either use a plain text editor, or be careful to save files as plain text.

What if we want to select lines from an arbitrary file? We could edit middle.sh each time to change the filename, but that would probably take longer than just retyping the command. Instead, let's edit middle.sh and make it more versatile:

```
$ nano middle.sh
```

Now, within “nano”, replace the text octane.pdb with the special variable called \$1:

```
head -n 15 "$1" | tail -n 5
```

Inside a shell script, \$1 means “the first filename (or other parameter) on the command line”. We can now run our script like this:

```
$ bash middle.sh octane.pdb
```

ATOM	9	H	1	-4.502	0.681	0.785	1.00	0.00
ATOM	10	H	1	-5.254	-0.243	-0.537	1.00	0.00
ATOM	11	H	1	-4.357	1.252	-0.895	1.00	0.00
ATOM	12	H	1	-3.009	-0.741	-1.467	1.00	0.00
ATOM	13	H	1	-3.172	-1.337	0.206	1.00	0.00

or on a different file like this:

```
$ bash middle.sh pentane.pdb
```

ATOM	9	H	1	1.324	0.350	-1.332	1.00	0.00
ATOM	10	H	1	1.271	1.378	0.122	1.00	0.00
ATOM	11	H	1	-0.074	-0.384	1.288	1.00	0.00
ATOM	12	H	1	-0.048	-1.362	-0.205	1.00	0.00
ATOM	13	H	1	-1.183	0.500	-1.412	1.00	0.00

✦ Double-Quotes Around Arguments

For the same reason that we put the loop variable inside double-quotes, in case the filename happens to contain any spaces, we surround \$1 with double-quotes.

We still need to edit `middle.sh` each time we want to adjust the range of lines, though. Let's fix that by using the special variables \$2 and \$3 for the number of lines to be passed to `head` and `tail` respectively:

```
$ nano middle.sh
```

```
head -n "$2" "$1" | tail -n "$3"
```

We can now run:

```
$ bash middle.sh pentane.pdb 15 5
```

```
ATOM      9  H          1      1.324   0.350  -1.332   1.00   0.00
ATOM     10  H          1      1.271   1.378   0.122   1.00   0.00
ATOM     11  H          1     -0.074  -0.384   1.288   1.00   0.00
ATOM     12  H          1     -0.048  -1.362  -0.205   1.00   0.00
ATOM     13  H          1     -1.183   0.500  -1.412   1.00   0.00
```

By changing the arguments to our command we can change our script's behaviour:

```
$ bash middle.sh pentane.pdb 20 5
```

```
ATOM     14  H          1     -1.259   1.420   0.112   1.00   0.00
ATOM     15  H          1     -2.608  -0.407   1.130   1.00   0.00
ATOM     16  H          1     -2.540  -1.303  -0.404   1.00   0.00
ATOM     17  H          1     -3.393   0.254  -0.321   1.00   0.00
TER       18          1
```

This works, but it may take the next person who reads `middle.sh` a moment to figure out what it does. We can improve our script by adding some **comments** at the top:

```
$ nano middle.sh
```

```
# Select lines from the middle of a file.
# Usage: bash middle.sh filename end_line num_lines
head -n "$2" "$1" | tail -n "$3"
```

A comment starts with a # character and runs to the end of the line. The computer ignores comments, but they're invaluable for helping people understand and use scripts.

What if we want to process many files in a single pipeline? For example, if we want to sort our `.pdb` files by length, we would type:

```
$ wc -l *.pdb | sort -n
```

because `wc -l` lists the number of lines in the files (recall that `wc` stands for ‘word count’, adding the `-l` flag means ‘count lines’ instead) and `sort -n` sorts things numerically. We could put this in a file, but then it would only ever sort a list of `.pdb` files in the current directory. If we want to be able to get a sorted list of other kinds of files, we need a way to get all those names into the script. We can’t use `$1`, `$2`, and so on because we don’t know how many files there are. Instead, we use the special variable `$@`, which means, “All of the command-line parameters to the shell script.” We also should put `$@` inside double-quotes to handle the case of parameters containing spaces (`"$@"` is equivalent to `"$1" "$2" ...`) Here’s an example:

```
$ nano sorted.sh
```

```
# Sort filenames by their length.
# Usage: bash sorted.sh one_or_more_filenames
wc -l "$@" | sort -n
```

```
$ bash sorted.sh *.pdb ../creatures/*.dat
```

```
9 methane.pdb
12 ethane.pdb
15 propane.pdb
20 cubane.pdb
21 pentane.pdb
30 octane.pdb
163 ../creatures/basilisk.dat
163 ../creatures/unicorn.dat
```

✦ Why Isn’t It Doing Anything?

What happens if a script is supposed to process a bunch of files, but we don’t give it any filenames? For example, what if we type:

```
$ bash sorted.sh
```

but don’t say `*.dat` (or anything else)? In this case, `$@` expands to nothing at all, so the pipeline inside the script is effectively:

```
$ wc -l | sort -n
```

Since it doesn’t have any filenames, `wc` assumes it is supposed to process standard input, so it just sits there and waits for us to give it some data interactively. From the outside, though, all we see is it sitting there: the script doesn’t appear to do anything.

We have two more things to do before we’re finished with our simple shell scripts. If you look at a script like:

```
$ wc -l "$@" | sort -n
```

you can probably puzzle out what it does. On the other hand, if you look at this script:

```
# List files sorted by number of lines.
$ wc -l "$@" | sort -n
```

you don't have to puzzle it out — the comment at the top tells you what it does. A line or two of documentation like this make it much easier for other people (including your future self) to re-use your work. The only caveat is that each time you modify the script, you should check that the comment is still accurate: an explanation that sends the reader in the wrong direction is worse than none at all.

Second, suppose we have just run a series of commands that did something useful — for example, that created a graph we'd like to use in a paper. We'd like to be able to re-create the graph later if we need to, so we want to save the commands in a file. Instead of typing them in again (and potentially getting them wrong) we can do this:

```
$ history | tail -n 5 > redo-figure-3.sh
```

The file `redo-figure-3.sh` now contains:

```
297 bash goostats -r NENE01729B.txt stats-NENE01729B.txt
298 bash goodiff stats-NENE01729B.txt /data/validated/01729.txt > 01729-differences.txt
299 cut -d ',' -f 2-3 01729-differences.txt > 01729-time-series.txt
300 ygraph --format scatter --color bw --borders none 01729-time-series.txt figure-3.png
301 history | tail -n 5 > redo-figure-3.sh
```

After a moment's work in an editor to remove the serial numbers on the commands, and to remove the final line where we called the `history` command, we have a completely accurate record of how we created that figure.

In practice, most people develop shell scripts by running commands at the shell prompt a few times to make sure they're doing the right thing, then saving them in a file for re-use. This style of work allows people to recycle what they discover about their data and their workflow with one call to `history` and a bit of editing to clean up the output and save it as a shell script.

Nelle's Pipeline: Creating a Script

An off-hand comment from her supervisor has made Nelle realize that she should have provided a couple of extra parameters to `goostats` when she processed her files. This might have been a disaster if she had done all the analysis by hand, but thanks to `for` loops, it will only take a couple of hours to re-do.

But experience has taught her that if something needs to be done twice, it will probably need to be done a third or fourth time as well. She runs the editor and writes the following:

```
# Calculate reduced stats for data files at J = 100 c/bp.
for datafile in "$@"
do
    echo $datafile
    bash goostats -J 100 -r $datafile stats-$datafile
done
```

(The parameters `-J 100` and `-r` are the ones her supervisor said she should have used.) She saves this in a file called `do-stats.sh` so that she can now re-do the first stage of her analysis by typing:

```
$ bash do-stats.sh *[AB].txt
```

She can also do this:

```
$ bash do-stats.sh *[AB].txt | wc -l
```

so that the output is just the number of files processed rather than the names of the files that were processed.

One thing to note about Nelle's script is that it lets the person running it decide what files to process. She could have written it as:

```
# Calculate reduced stats for A and Site B data files at J = 100 c/bp.
for datafile in *[AB].txt
do
    echo $datafile
    bash goostats -J 100 -r $datafile stats-$datafile
done
```

The advantage is that this always selects the right files: she doesn't have to remember to exclude the 'Z' files. The disadvantage is that it *always* selects just those files — she can't run it on all files (including the 'Z' files), or on the 'G' or 'H' files her colleagues in Antarctica are producing, without editing the script. If she wanted to be more adventurous, she could modify her script to check for command-line parameters, and use *[AB].txt if none were provided. Of course, this introduces another tradeoff between flexibility and complexity.

 Variables in Shell Scripts 

 List Unique Species 

 Find the Longest File With a Given Extension 

 Why Record Commands in the History Before Running Them? 

 Script Reading Comprehension 

 Debugging Scripts 

Key Points

- Save commands in files (usually called shell scripts) for re-use.
- `bash filename` runs the commands saved in a file.
- `$@` refers to all of a shell script's command-line parameters.
- `$1`, `$2`, etc., refer to the first command-line parameter, the second command-line parameter, etc.
- Place variables in quotes if the values might have spaces in them.
- Letting users decide what files to process is more flexible and more consistent with built-in Unix commands.

Copyright © 2016 Software Carpentry Foundation (<https://software-carpentry.org>)

Source (<https://github.com/swcarpentry/shell-novice/>) / Contributing
(<https://github.com/swcarpentry/shell-novice/blob/gh-pages/CONTRIBUTING.md>) /
Contact (<mailto:lessons@software-carpentry.org>)