

Intro to Shell Scripting

Dr. Princess Rodriguez

2025-01-29

Learning Objectives

- Capture commands into a shell script
- Implement variables in a shell script

Recap from last class

- `wc` counts lines, words, and characters in its inputs
- `command > [file]` redirects a command's output to a file (overwriting any existing content)
- `command >> [file]` appends a command's output to a file
- `[first] | [second]` is a pipeline: the output of the first command is used as the input to the second

cut

- **cut is a command that extracts columns from files.**

We can use `cut` with the `-f` argument to specify which specific fields or columns from the dataset we want to extract. Let's say we want to get the 1st and 3rd column of a file, we can use:

```
cut -f 1,3 file.csv
```

Please note: the argument `-f 1,3` cannot be written with spaces between the field numbers. If you write it as `-f 1, 3` (with a space), the `cut` command will treat 1, as one argument and 3 as another, leading to an error.

```
cut -f 1, 3 file.csv
```

```
# Error: "cut: invalid byte, character, or field list"
```

“Class Exercise: Participation Grade”

Before moving on, please complete the following class act

[Class-activity] (<https://forms.gle/z16AAU9oxNg9Dohk6>)

sort

- **sort** is a command used to sort lines in text files in a particular order.

sort options

Option	Description	Example Usage
-u (unique)	Outputs unique lines	<code>sort -u [FILE]</code>
-k (key)	Allows sorting based on a specific key field/column in each line	<code>sort -k 2 [FILE]</code>
-n (numerics)	Performs a numeric sort	<code>sort -n [FILE]</code>

What is a BED File?

A BED file is a widely used format in genomics to store genomics regions and associated annotations. It is primarily used for specifying intervals or features on a reference genome such as genes, exons, or regulatory elements. A BED file consists of at least three required columns but can include additional optional columns.

Required Columns for BED:

- 1 Chromosome: The name of the chromosome (chr1)
- 2 Start: The start position of a feature
- 3 End: The end position of a feature

chr1	213941196	213942363
chr1	213942363	213943530
chr1	213943530	213944697

chr7	127471196	127472363	Pos1	0	+
chr7	127472363	127473530	Pos2	0	+
chr7	127473530	127474697	Pos3	0	+

Chromosome ID →

Start location ↑ End location ↑ Name ↑ Strand ↑

Phase (reading frame)

Figure 1: Bed File

Class Exercise #2

Dr. Patel: “Emma I noticed that the file `coordinates_PLEKHN1.txt` contained a number of duplicates. Can you please modify this file and save it as a BED file so I can upload it to UCSC Genome Browser. The file extension is `.bed`. Thanks.”

Shell scripts

Over the past few weeks, you have been introduced to a number of commands to explore data files. To demonstrate the function of each command we have run them one at a time at the command prompt. The command prompt is useful for testing out commands and also performing simple tasks like exploring and organizing the file system. However, when we are running bioinformatic analyses which require a series of tasks to be run, there is a more efficient way to do this.

Shell scripts continued

Shell scripts are **text files that contain commands we know we want to run.**

What can shell scripts be used for?

- for creating, maintaining and implementing system-wide scripts
- automating tedious repetitive tasks
- scheduling and executing system tasks
- for automating the installation process for new software or for new software updates across the organization
- for scheduling data backup process

Writing a simple shell script

We are ready to see what makes the shell such a powerful programming environment. To create our first script, we are going to take some of the commands we have run previously and save them into a file so that we can **re-run all those operations** again later, by typing just **one single command**.

Writing a simple shell script

When working on the command line you can give files any extension (.txt, .tsv, .csv, etc.). Similarly, for a shell script you don't need a specific extension. However, it is **best practice to give shell scripts the extension .sh** (bash shell script file).

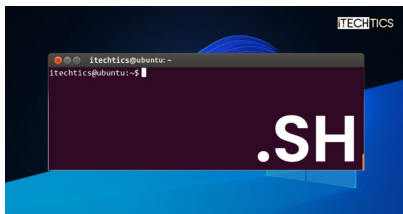


Figure 2: SH extension

“Class Exercise #3

You will have ~5 minutes to complete. NOTE: It is never my intention to rush you. If you do find you need more time to participate, please let me know.**

Objective: Create a script that contains 3 lines.

Final Question: *Dr.Patel: “Emma thanks for carrying out that exercise. How many FASTQ reads in Mov10 oe_1.subset.fq contained 10 consecutive N’s or more?”*

****Submit your answer in Class Participation survey.****

One FASTQ Read

```
<figure markdown="span">  
  ![One FASTQ Read](../img/fastq_fig.jpg){ width="200"}  
  <figcaption> One FASTQ Read </figcaption>  
</figure>
```


Executing Scripts #1

There are two main ways to execute a script:

1. Direct Invocation by the Shell: When you use `sh script-name.sh`, the shell (*i.e. way to interact with kernel*) reads and executes the script file. The executable permissions of the script file is not checked because we are directly passing it as an argument to the shell interpreter.

- `sh script.sh`: The script does not need executable permissions as `sh` is invoked manually
- Best for testing scripts

Executing Scripts #2

2. Executable Permissions: Making a script executable allows you to run it directly using `./script-name.sh`. Here, the kernel (*core of the OS, manages hardware and system resources*) will check if the script file has executable (x) permissions and invokes the interpreter specified in the scripts shebang line.

- `./script.sh`: Requires the script file to have executable permissions and the kernel uses the interpreter specified in the shebang line.
- Best for automation and reusable scripts

Bash VARIABLE

A **variable** is a common concept shared by many programming languages. **Think of variables as a temporary store or *bucket* for a piece of information.** This bucket will have a name associated with it therefore when referring to the information inside the bucket, we can use the name of the bucket instead!

Creating a Bash VARIABLE

TO create a variable in bash, you will provide the name of the variable, followed by the equals sign and finish with the value we want to assign to the variable.

```
name_of_variable=value_of_variable
```

Important notes when creating Bash VARIABLES

- Note that the variable name cannot contain spaces, nor can there be spaces on either side of the equals sign.
- The variable name can have only letters (a to z or A to Z), numbers (0 to 9), or the underscore character (_). The wrong character usage in the variable name will cause a syntax error.
- By convention, the variable names in UNIX are in UPPERCASE.

Class Exercise

Create a variable called NUM that has the number 25 stored inside it:

```
NUM=25
```

If we are using our bucket analogy - You can think of the variable NUM like this:



What happens after we created NUM?

Once you press return, you will find yourself back at the command prompt. But nothing happened. . . so how do we know that we actually created a variable?

One way to see the variable created is by using the `echo` command.

echo command

The echo command is used to display text or the value of variables to the terminal/standard output.

```
echo "Hello, World!"
```

```
## Hello, World!"
```


echo command with variables

```
NAME=ALICE  
echo "Hello, $NAME!"
```

```
## Hello, Alice!
```

To display the contents of the variable we need to **explicitly use a \$ in front of the variable name**:

echo with NUM

```
echo $NUM
```

You should see the number 25 returned to you. Notice that when we created the variable, we did not use the \$.

What is the \$

The \$ is a standard shell notation for defining and using variables. The \$ tells the shell interpreter to treat the variable as a variable name and substitute its value in its place, rather than treat it as text or an external command.

Therefore, when defining a variable (i.e. setting the value) you can just type it as is, but when retrieving the value of a variable you must use the \$!

Important note when dealing with variables

NOTE: Variables are not physical entities like files. When you create files you can use `ls` to list contents and see if the file exists. To list all variables in your environment you can use the command `declare` with the `-p` option. You will notice that while you only have created one variable so far, the output of `declare -p` will be more than just one variable. These other variables are called environment variables. To remove a variable you can use `unset`.

Using variables as input to commands

One important aspect of the variable is that the value stored inside can be used as input to commands.

Let's solidify this important concept with Class Exercise #4.

Class Exercise #4

- 1 Create a new variable called `FILE`. Use the name of one of the fastq files in the `raw_fastq` directory as the value of the variable.
- 2 Recall the variable with `echo`
- 3 Check the number of lines in the `FILE` variable.

Another important note when dealing with variables

NOTE: The variables we create in a session are system-wide, and independent of where you are in the filesystem. This is why we can reference it from any directory. However, it is only available for your current session. If you exit the cluster and login again at a later time, the variables you have created will no longer exist.

Utility of variables

- When creating shell scripts, variables are used to store information that can be used later in the script (once or many times over).
- The value stored can be hard-coded in as we have done above, assigning the variable a numeric or character value.
- Alternatively, the value stored can be the output of another command.

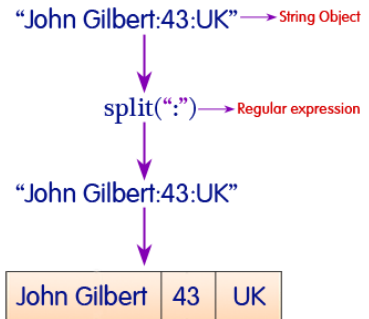
We will demonstrate this using a new command called `basename`.

basename

The **basename command** is used to extract the file name or directory name from a given file path. This is accomplished using **string splitting**.

String splitting

String splitting is a way to break a larger string into smaller parts based on a specified delimiter.



© w3resource.com

Figure 4: String Splitting

Other “common” strings used include:

- Space: Useful for splitting words in a sentence.
- Tab (`\t`): Often used in tab-delimited data files.
- Semicolon (`;`): Another popular choice for separating values in data.
- Colon (`:`): Commonly used in key-value pairs.
- Pipe (`|`): Used in various data formats, such as CSV files.
- Hyphen (`-`): Can be used to split ranges or parts of a string.
- Underscore (`_`): Frequently used in variable or function names.
- Forward dash (`/`): Useful for splitting file paths.

Basic Usage of basename:

1. Extract file name

```
basename /path/to/file.txt
```

```
## file.txt
```

Class Example with basename

```
basename ~/unit1_unix/raw_fastq/Mov10_oe_1.subset.fq
```

The command will only return the file name.

2. Remove file extension

```
basename /path/to/file.txt .txt
```

```
## file
```

“Class Exercise #5”

Use `basename` with the file `Irrel_kd_1.subset.fq` as input.
Return only `Irrel_kd_1` to the terminal.

Storing the `basename` output in a variable

The `basename` command returns a character string and this too can be stored inside a variable. To do this without error, we need to add another special syntax because when we run the command we will generate spaces. If you remember earlier, one of the rules of creating variables is that there cannot be any spaces.

NOTE: The backtick key `. On most keyboards this character is located just underneath the `esc` key.

Storing the `basename` output in a variable

example

```
VARIABLE=`basename /path/to/file`
```

Let's try an example:

```
samplename=`basename ~/unit1_unix/raw_fastq/Mov10_oe_1.su
```

Check to see what got stored in the `samplename` variable:

```
echo $samplename
```

Utility of `basename` command

- It is hard to see the utility of this command by just running it at command-line, but it is very useful command when creating scripts for analysis.
- Within a script it is common to create an output file and the `basename` allows us to easily create a prefix to use for naming the output files.

Shell scripting using Jupyter Notebook

Now it's time to put all of these concepts together to create a more advanced version of the script. This script will allow the user to get information on any given directory. These are the steps you will code into a shell script using Jupyter Notebook:

- ➊ Assign the path of the directory to a variable
- ➋ Create a variable that stores only the directory name (and no path information)
- ➌ Move from the current location in the filesystem into the directory we selected in 1.
- ➍ List the contents of the directory
- ➎ List the total number of files in the directory

It seems like a lot, but you are equipped with all the necessary concepts and commands to do this quite easily!

What is Jupyter Notebook

Jupyter Notebook is an open-source web application that allows you to create and share documents containing live code, equations, visualizations, and narrative text. It's widely used in data science, scientific research, and education. The term “Jupyter” is derived from the combination of three core programming languages it supports: Julia, Python, and R.



Figure 5: Jupyter Logo

Request time to use Jupyter Notebook on the VACC

- Slurm Account: mmg3320
- Partition: general
- Everything else leave as default
- Press Launch

This might take a few minutes to start up!

Class Exercise #6 and your Homework Assignment Part B

This is a self-paced assignment. This is also the final assignment for today. If you would like to do this from home feel free! You will need to submit (2) screenshots from this exercise with your homework for this week.

- 1 To get started move into the other directory. Jupyter Notebook is user-friendly. You should be able to click from `unit1_unix` into other easily.
- 2 Press New (Right-side) -> New file -> and create a script called `directory_info.sh`.

Summary

In today's lesson, we described shell scripts and introduced a few related concepts that are helpful when you are starting out. It is important to understand each of the individual concepts, but also to see how they all come together to add flexibility and efficiency to your script. Later on we will further illustrate the power of scripts and how they can make our lives (when coding) much easier. Any type of data you will want to analyze will inevitably involve not just one step, but many steps and perhaps many different tools/software programs. Compiling these into a shell script is the first step in creating your analysis workflow!

Citation

This lesson has been developed by members of the teaching team at the Harvard Chan Bioinformatics Core (HBC). These are open access materials distributed under the terms of the Creative Commons Attribution license (CC BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

- *The materials used in this lesson were derived from work that is Copyright © Data Carpentry (<http://datacarpentry.org/>). All Data Carpentry instructional material is made available under the Creative Commons Attribution license (CC BY 4.0).*
- *Adapted from the lesson by Tracy Teal. Original contributors: Paul Wilson, Milad Fatenejad, Sasha Wood and Radhika Khetani for Software Carpentry (<http://software-carpentry.org/>)*
- *Other Authors: Meeta Mistry, Bob Freeman, Mary Piper, Radhika Khetani, Jie Liu, Will Gammerding*