

1 Linux for Bioinformatics

1.1 Introduction

Unix and Unix-like operating systems are the standard operating system on most large computer systems in scientific research, in the same way that Microsoft Windows is the dominant operating system on desktop PCs. In this course we will use Linux, a Unix-like operating system which was originally created to provide a free open source operating system for PCs.

Linux is a powerful, secure, robust and stable operating system which allows dozens of people to run programs on the same computer at the same time. This is why it is the preferred operating system for large-scale scientific computing. It runs on all kinds of machines, from mobile phones (Android), desktop PCs... to supercomputers.

1.1.1 Why Linux?

Increasingly, the output of biological research exists as *in silico* data, usually in the form of large text files. Linux is particularly suitable for working with such files and has several powerful and flexible commands that can be used to process and analyse this data. One advantage of learning Linux is that many of the commands can be combined in an almost unlimited fashion. So if you can learn just six Linux commands, you will be able to do a lot more than just six things.

Linux contains hundreds of commands, but to conduct your analysis you will probably only need 10 or so to achieve most of what you want to do. In this tutorial we will introduce you to some useful Linux commands and provide examples of how they can be used in bioinformatics analyses.

1.2 Learning outcomes

By the end of the tutorial you can expect to be able to:

- Understand the Linux directory structure and navigate around this structure
- Extract information from large files
- Use regular expressions to search for particular patterns in a file
- Create a bash script to perform several tasks at once

1.3 Tutorial sections

This tutorial comprises the following sections:

1. [Basic Linux](#)
2. [Commands grep and awk](#)
3. [Advanced Linux \(loops and Bash scripts\)](#)
4. [Bash scripting](#)

Note: If you are an experienced linux user then please feel free to skip the [Basic Linux](#) section and move directly to the section [Commands grep and awk](#).

Also, do not worry if you do not complete all sections in the time allocated, a good target to aim for is the end of section on [Advanced Linux \(loops and Bash scripts\)](#). The remaining sections are optional and are for students who would like to expand their Linux skills and can be completed outside the course hours.

1.4 Authors and License

This tutorial was created by [Jacqui Keane](#) and [Martin Hunt](#).

The content is licensed under a [Creative Commons Attribution 4.0 International License \(CC-BY 4.0\)](#).

1.5 Running the commands in this tutorial

You can follow this tutorial by typing all the commands you see in a terminal window on your computer. This is similar to the “Command Prompt” window on MS Windows systems, which allows the user to type DOS commands to manage files.

To get started, open a terminal window and type the command below followed by the Enter key:



```
cd ~/course_data/linux/data
```

Now you can follow the instructions in the tutorial from here.

1.6 Cheat sheet

We’ve also included a [cheat sheet](#) at the end of this tutorial. It probably won’t make a lot of sense now, but it might be a useful reminder later.

1.7 Let’s get started!

To get started with the tutorial, go to the next section: [Basic Linux](#)

2 Basic Linux

2.1 The Commandline

The commandline or ‘terminal’ is a text-based interface you can use to run programs and analyse your data. If this is your first time using one it will seem pretty daunting at first but, with just a few commands, you’ll start to see how it helps you to get things done more efficiently.

2.2 Getting started

When working on the commandline you will always work in a specific location or directory in the filesystem. So let’s check that you’re in the right place. Type the command below in the terminal window followed by the Enter key:



```
pwd
```

The `pwd` command stands for print working directory. It should display something similar to:

```
/home/manager/course_data/linux/data
```

This means you are working in the `data` directory found inside the `linux` directory inside the `course_data` directory and so on. There will be more on the linux filesystem structure shortly.

Now continue through this tutorial, entering any commands that you encounter (highlighted in a grey box with a keyboard symbol) into your terminal window followed by the Enter key. Let’s start by listing the contents of the current working directory:



```
ls
```

You should see 6 items:

```
advanced_linux  awk  bash_scripts  basic  files  grep
```

2.3 Commands, options and arguments

A command consists of three parts: the *command*, *options*, and *arguments*. Of these, the command is required while options and arguments are optional. Options are used to modify the default behavior of the command. Arguments are used to provide input to the command.

```
command -options arguments
```

In the previous call to `ls`, the command is `ls` and there are no options or arguments passed to the `ls` command.

Now let’s add some options, here the command is `ls`, and `-a` and `-l` are options:



```
ls -a -l
```

Note that there is a space between the command `ls` and the options `-a` and `-l`. There is no space between the dashes and the letters `a` and `l`.

The `-a` option will list *all* contents of a directory including hidden files and directories (hidden files and directories are not shown by default and are used to help important data from being deleted) and the `-l` option will print additional information about each item in the directory (size, owner, date changed etc.).

You can also combine the `-a` and `-l` options into what appears to be a single `-al` option. It's almost always ok to do this for options which are made up of a single dash followed by a single letter.



```
ls -al
```

Now let's add an argument, here the command is `ls`, the options are `-a` and `-l` and `Pfalciparum` is the argument:



```
ls -al Pfalciparum/
```

The argument instructs the `ls` command to display the content of the specified directory, in this case `Pfalciparum` instead of the current directory.

There are some general points to remember that will make your life easier:

- Linux is case sensitive - typing `ls` is not the same as typing `LS`.
- `ls` is the letter `l` followed by the letter `s` (not the number one).
- Often when you have problems with Linux, it is due to a spelling mistake. Check that you have not misspelled a command or missed a space between the command and the options and arguments. Pay careful attention if typing commands across multiple lines.

2.4 Files and directories

Directories are the Linux equivalent of folders on a PC or Mac. They are organised in a hierarchy, so directories can have sub-directories and so on. Directories are very useful for organising your work and keeping your account tidy - for example, if you have more than one project, you can organise the files for each project into different directories to keep them separate.

The start of the filesystem is known as the root directory and denoted with `/` symbol. The location or directory that you are in at any given time is referred to as the current working directory denoted with the `.` symbol.

Every file and directory on the computer is found in a specific location. The location is specified as a path to that file or directory and can be expressed as either the *relative path* or *absolute path*.

The *relative path* is the location of a file or directory from the current working directory (`.`).

The *absolute path* is the location of a file or directory from the start of the file system or root directory (`/`).

For example, for the file called `directory_structure.png` under the `linux` directory, the location or *relative path* can be expressed as:

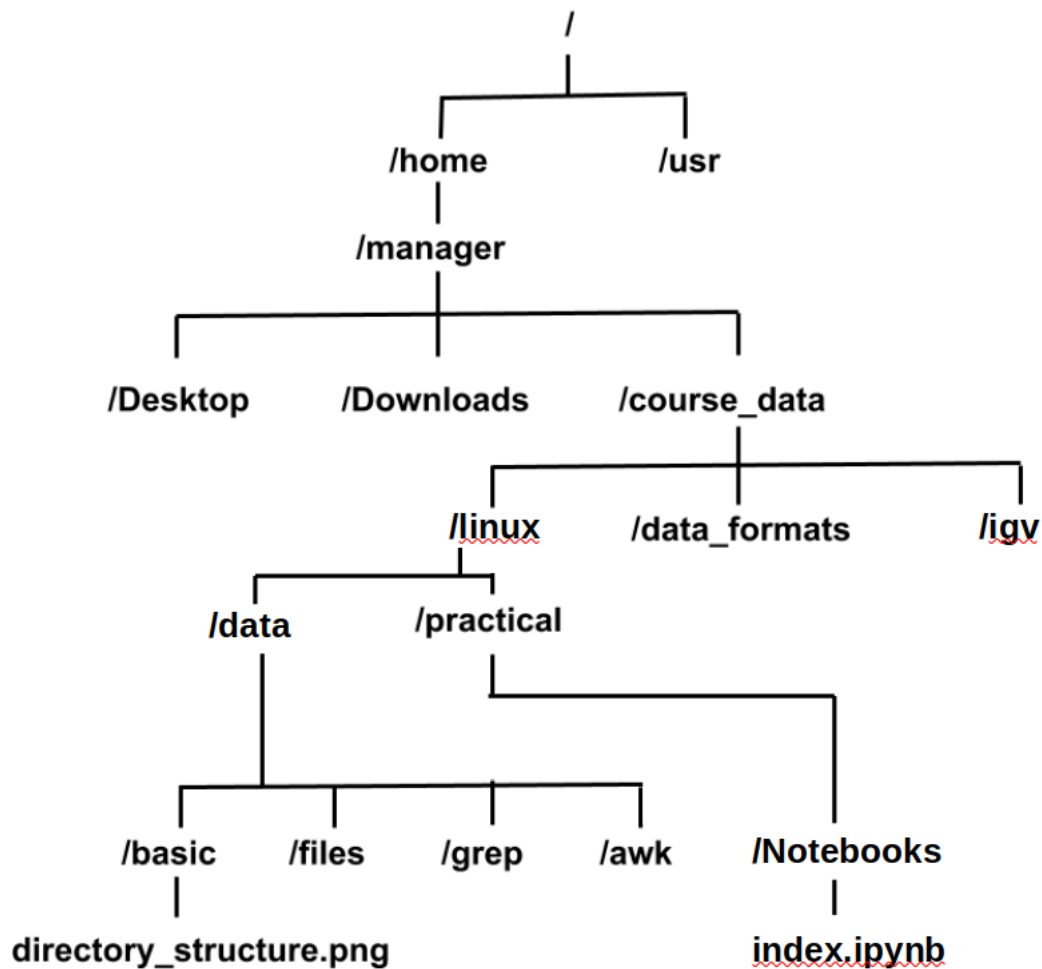


```
ls ../basic/directory_structure.png
```

The location of this file can also be expressed using the *absolute path* as:



```
ls /home/manager/course_data/linux/data/basic/directory_structure.png
```



2.5 tree - display the directory hierarchy

The command `tree` can be used to recursively list or display the content of a directory in a tree-like format. It outputs the directory paths and files in each sub-directory and a summary of a total number of sub-directories and files.

To display the contents of the current directory, type:



```
tree
```

To display the hierarchy of the 'linux' directory:




```
tree /home/manager/course_data/linux
```

2.6 cd - change directory

The command `cd` stands for change directory. The `cd` command will move you from the current working directory to another directory.

To move into the `basic` directory type the following command. Note, you'll remember this more easily if you type this rather than copying and pasting.


 `cd basic`

Now look at the directory hierarchy of this directory:

 `tree`

2.7 Getting help man

It is not possible to remember all the options and arguments for all Linux commands! A useful command to obtain further information on any Linux command is the `man` command. For example, to get a full description and examples of how to use the `find` command type the following command in a terminal window.


 `man find`

To exit out of the `man` command type `q`.


2.8 find - find a file

The `find` command can be used to find files in a directory hierarchy that match a specific criteria. For example, it can be used to recursively search the directory tree for a specific file name or pattern, seeking files and directories that match the given name or pattern.

To find all files in the current directory (`.`) named `Malaria.fa`:

 `find . -name "Malaria.fa"`

To find all files in the current directory (`.`) and all its subdirectories that have the suffix `.bed`:

 `find . -name "*.bed"`

How many `bed` files did you find?

Can you construct a command to find all the subdirectories contained in the current directory? **Hint:** You may need to go back to reading the manual page for the `find` command and search for the `-type` option. Don't spend too much time on this as the solution will be provided and explained later.



How many subdirectories did you find?

2.9 Tab completion

Typing out directory or file names is really boring and you're likely to make typos which will at best make your command fail with a strange error and at worst overwrite some of your carefully crafted analysis. *Tab completion* is a trick which normally reduces this risk significantly.

List the contents of the Styphi directory:



```
ls Styphi
```

Now instead of typing out `ls Styphi/`, try typing `ls S` and then press the tab key (instead of Enter). The rest of the folder name should just appear.

If you have two files or directories with similar names (e.g. `directory_structure.png` and `directory_structure2.png`) then you might need to give your terminal a bit of a hand to work out which one you want. Type the following and then press the tab key (instead of Enter).



```
ls -l d
```

In this case, when you press tab the terminal reads `ls -l directory_structure`, you then need to type 2 followed by another tab and it works out that you meant `directory_structure2.png`.

2.10 Tips

There are some short cuts for referring to directories:

- `.` Current directory (one full stop)
- `..` Directory above in the hierarchy (two full stops)
- `~` Home directory (tilde)
- `/` Root of the file system (like C: in Windows)

Try the following commands, what do they do?



```
ls .
```



```
ls ..
```



```
ls ~
```

2.11 Exercises

Many people panic when they are confronted with a Linux prompt! Don't! All the commands you need to solve these exercises are provided above and don't be afraid to make a mistake. If you get lost ask an instructor. If you are a person skilled at Linux, be patient this is only a short exercise.

To begin, open a terminal window and navigate to the `basic` directory under the `linux` directory (remember use the command `cd`) and then complete the exercise below.

1. List *all* the contents of the `basic` directory. How many files did you find?
2. How many files are there in the `Styphi` directory?

3. How many files are there in the `Pfalziparum` directory?
4. Use the `find` command to find all `gff` files in the `linux` directory, how many did you find?
5. Use the `find` command to find all `fasta` files in the `linux` directory, how many did you find?

When you have completed these exercises move on to the next part of the tutorial, [grep and awk](#).

3 The commands `grep` and `awk`

In this section we will introduce the commands `grep` and `awk`.

3.1 Searching inside files with `grep`

A common task is to extract information from large files. This can be achieved using the Linux command `grep`, which stands for “Globally search for a Regular Expression and Print”. The meaning of this acronym will become clear later, when we discuss Regular Expressions. First, we will consider simpler examples.


Before we start, change into the `grep` directory:



```
cd ../grep
```

3.1.1 Simple pattern matching

We will use a small example file (in “BED” format), which contains information about gene features in a genome. Take a look at the contents of the first example BED file used in this section:



```
cat gene_expression.bed
```

A bed file is a column-based file, each column contains information about the feature and is separated by a tab character. There can be more than 10 columns, but only the first three are required to be a valid file. The file format is described in full here: <http://genome.ucsc.edu/FAQ/FAQformat#format1>. We will use the first 5 columns:

1. Sequence name
2. start position (starting from 0, not 1)
3. end position (starting from 0, not 1)
4. feature name
5. score (which is used to store the gene expression level in our examples).


In reality, such a file could contain 100,000s of lines, so that it is not practical to read manually. Suppose we are interested in all the genes from chromosome 2. We can find all these lines using `grep`:



```
grep chr2 gene_expression.bed
```

This has shown us all the lines that contain the text or string “chr2”.


Now we want to find all genes from chromosome 2 that are on the positive strand. We can use the `|` symbol which is also known as the pipe symbol to extract the genes that are on the positive strand, using `grep` a second time:



```
grep chr2 gene_expression.bed | grep +
```


This *pipes* (sends) the output of `grep chr2 gene_expression.bed` into the input of `grep +`.

Since `grep` is reporting a match to a string *anywhere* on a line, such simple searches can have undesired consequences. For example, consider the result of doing a similar search for all the genes in chromosome 1:


```
 grep chr1 gene_expression.bed
```

Oops! We found genes in chromosome 10, because “chr1” is a substring (subset) of “chr10”.

Or consider the following file, where the genes have unpredictable names (which is not unusual for bioinformatics data).

```
 cat gene_expression_sneaky.bed
```

Now we try to find genes on chromosome 1 that are on the negative strand. We put the minus sign in quotes, to stop Linux interpreting this as an option to `grep`, as opposed to the string we are searching for:

```
 grep chr1 gene_expression_sneaky.bed | grep '-'
```


The extra lines are found by `grep` because of matches in columns we were not expecting to match. Remember, `grep` is reporting these lines because they each contain the strings “chr1” and “-” *some-where*.


We need a way to make searching with `grep` more specific.

3.1.2 Regular expressions


Regular expressions provide the solution to the above problems. They are ways of defining more specific patterns to search for.

Matching the start and end of lines First, we can specify that a match must be at the start of a line using the symbol “^”, which means “start of line”. Without the ^, we find any match to “chr1”:

```
 grep chr1 gene_expression_sneaky.bed
```

```
 grep '^chr1' gene_expression_sneaky.bed
```

Good! We have removed the match to the badly-named gene “chr11.gene1”, which is on chromosome 8. Now we want to avoid matching chromosomes 10 and 11. This can be done by also looking for a “tab” character, which is represented by writing `\t`. For technical reasons, which are beyond the scope of this course, we must also put a dollar sign before the quotes to make any search involving a tab character work.

```
 grep '$'^chr1\t' gene_expression_sneaky.bed
```

To find the genes on the negative strand, all that remains is to match a minus sign at the *end* of the line (so that we do not find “sneaky-gene3”). We can do this using the dollar “\$”, which means “end of line”.

```
grep $'^chr1\t' gene_expression_sneaky.bed | grep '\-$'
```

Wildcards and alphabets Another special character in regular expressions is the dot: “.”. This stands for any single character. For example, this finds all matches to chromosomes 1-9, and chromosomes X and Y:

```
grep $'^chr.\t' gene_expression.bed
```

In fact, the earlier command that found all genes on chromosome 1 that are on the negative strand, could be found with a single call to `grep` instead of two calls piped together. To do this, we need a regular expression that finds lines that:

- start with `chr1`, then a tab character
- end with a minus
- have arbitrary characters between.

The asterisk “*” has a special meaning: it says to match any number (including zero) of whatever character is before the *. For example, the regular expression ‘AC*G’ will match AG, ACG, ACCG, etc. The simpler, improved command is:

```
grep $'^chr1\t.*-$' gene_expression_sneaky.bed
```

As well as matching any character using a dot, we can define any list of characters to match, using square brackets. For example, `[12X]` means match a 1, 2, or an X. This can be used to find all genes from chromosomes 1, 2 and X:

```
grep $'^chr[12X]\t' gene_expression.bed
```

Or just the autosomes may be of interest. To do this we introduce two new features:

- Ranges can be given in square brackets, for example `[1-5]` will match 1, 2, 3, 4 or 5.
- The plus sign “+” has a special meaning that is similar to “*”. Instead of any number of matches (including zero), it looks for at least one match. To avoid simply matching a plus sign, it must be preceded by a backslash: “\+”. For example, the regular expression ‘AC\+G’ will match ACG, ACCG, ACCCG etc (but will not match AG).


Warning: Adding a backslash is often called *escaping* (e.g. *escape the plus symbol*). Depending on the software you’re using (and the options you give it), you may need to escape the symbol to indicate that you want its special regex meaning (e.g. multiple copies of the last character please) or its literal meaning (e.g. give me a ‘+’ symbol please). If your command isn’t working as you expect, try playing with these options and always test your regular expression before assuming it gave you the right answer.

The command to find the autosomes is:


```
grep $'^chr[0-9]\+\t' gene_expression.bed
```

3.1.3 Other `grep` options


The Linux command `grep` and regular expressions are extremely powerful and we have only scratched the surface of what they can do. Take a look at the manual (by typing `man grep`) to get an idea. A few particularly useful options are discussed below.

 `man grep`


Counting matches A common scenario is counting matches within files. Instead of output each matching line, the option “`-c`” tells `grep` to report the number of lines that matched. For example, the number of genes in the autosomes in the above example can be found by simply adding `-c` to the command.

 `grep -c $'^chr[0-9]\+\t' gene_expression.bed`

Case sensitivity By default, `grep` is case-sensitive. It can be useful to ignore the distinction between upper and lower case using the option “`-i`”. Suppose we have a file of sequences, and want to find the sequences that contain the string `ACGT`. It is not unusual to come across files that have a mix of upper and lower case nucleotides. Consider this FASTA file:

 `cat sequences.fasta`


A simple search for `ACGT` will not return all the results:

 `grep ACGT sequences.fasta`

However, making the search case-insensitive solves the problem.

 `grep -i ACGT sequences.fasta`

Inverting matches By default, `grep` reports all lines that do match the regular expression. Sometimes it is useful to filter a file, by reporting lines that *do not* match the regular expression. Using the option “`-v`” makes `grep` “invert” the output. For example, we could exclude genes from autosomes in the BED file from earlier.


 `grep -v $'^chr[0-9]\+\t' gene_expression.bed`

3.1.4 Replacing matches to regular expressions

Finally, we show how to replace every match to a regular expression with something else, using the command “`sed`”. The general form of this is:

```
sed 's/regular expression/new string/' input_file
```

This will output a new version of the input file, with each match to the regular expression replaced with “`new string`”. For example, try:



```
sed 's/^chr/chromosome/' gene_expression.bed
```

This will replace every occurrence of the text “chr” that appears at the start of a line in the file `gene_expression.bed` with the text “chromosome”.

3.1.5 Exercises

The following exercises all use the FASTA file `exercises.fasta`. Before starting the exercises, open a new terminal and navigate to the `linux/data/grep` directory, which contains `exercises.fasta`.

Use `grep` to find the answers. Hint: some questions require you to use `grep` twice, and possibly some other Linux commands.

1. Make a `grep` command that outputs just the lines with the sequence names.
2. How many sequences are in the file?
3. Do any sequence names have spaces in them? What are their names?
4. Make a `grep` command that outputs just the lines with the sequences, not the names.
5. How many sequences contain unknown bases (an “n” or “N”)?
6. Are there any sequences that contain non-nucleotides (something other than A, C, G, T or N)?
7. How many sequences contain the 5’ cut site GCWGC (where W can be an A or T) for the restriction enzyme *AceI*?
8. Are there any sequences that have the same name? You do not need to find the actual repeated names, just whether any names are repeated. (Hint: it may be easier to first discover how many unique names there are).

3.2 File processing with AWK

AWK is a programming language named after the initials of its three inventors: Alfred Aho, Peter Weinberger, and Brian Kernighan. AWK is incredibly powerful at processing files, particularly column-based files (e.g. BED, GFF and SAM files), which are commonplace in Bioinformatics.


Before we start, change into the `awk` directory:



```
cd ../awk
```

3.2.1 Extracting columns from files

The command `awk` reads a file line-by-line, splitting each line into columns. This makes it easy to do simple things like extract a column from a file. We will use the following GFF file for our examples.



```
cat genes.gff
```


A GFF file is similar to a BED file in that it contains information about gene features in a genome. It is a column-based file, each column contains information about the feature. The columns in the GFF file are separated by tabs and each column has the following meaning:

1. Sequence name
2. Source - the name of the program that made the feature
3. Feature - the type of feature, for example gene or CDS

4. Start position
5. Stop position
6. Score
7. Strand (+ or -)
8. Frame (0, 1, or 2)
9. Optional extra information, in the form `key1=value1;key2=value2;...`

The score, strand, and frame can be set to `.` if it is not relevant for that feature. The final column 9 may or may not be present and could contain any number of key, value pairs.


We can use `awk` to just print the first column of the file. `awk` calls the columns `$1`, `$2`, ... etc, and the complete line is called `$0`. Try

```
 awk -F"\t" '{print $1}' genes.gff
```

A little explanation is needed.


- The option `-F"\t"` was needed to tell `awk` that the columns are separated by tabs (more on this later).
- For each line of the file, `awk` does what is inside the curly brackets. In this case, we simply print the first column (`$1`).

The repeated chromosome names are not nice. It is more likely to want to know just the unique names, which can be found by piping into the Linux command `sort -u`.

```
 awk -F"\t" '{print $1}' genes.gff | sort -u
```

In this command the `|` symbol is known as the *pipe* symbol. This *pipes* (sends) the output of the `awk` command into the input of `sort -u`. The `sort` will sort the contents of the input. When you sort the input, lines with identical content end up next to each other in the output. The `-u` option of the `sort` command will select only the unique values in the output.

You can connect as many commands as you want. For example, type:

```
 awk -F"\t" '{print $1}' genes.gff | sort -u | wc -l
```

This will count the number of lines in the output.

3.2.2 Filtering the input file

Similarly to `grep`, `awk` can be used to filter out lines of a file. However, since `awk` is column-based, it makes it easy to filter based on properties of any columns of interest. The filtering criteria can be added before the braces. For example, the following extracts just chromosome 1 from the file.

```
 awk -F"\t" '$1=="chr1" {print $0}' genes.gff
```

There are two important things to note from the above command:

1. `$1=="chr1"` means that column 1 must be *exactly* equal to `"chr1"`. This means that `"chr10"` is not found.

2. The “`{print $0}`” part only happens when the first column is equal to “chr1”, otherwise `awk` does nothing (the line gets ignored).

Awk commands are made up of two parts, a *pattern* (e.g. `$1=="chr1"`) and an *action* (e.g. `print $0`) which is contained in curly braces. The *pattern* defines which lines the *action* is applied to.

In fact, the action (the part in curly braces) can be omitted in this example. `awk` assumes that you want to print the whole line, unless it is told otherwise. This gives a simple method of filtering based on columns.

3.2.3 Exercises

The following exercises all use the BED file `exercises.bed`. Before starting the exercises, open a new terminal and navigate to the `awk` directory, which contains `exercises.bed`.

Use `awk` to find the answers to the following questions about the file `exercises.bed`. Many questions will require using pipes (eg “`awk ... | sort -u`” for question 1).

1. What are the names of the contigs in the file?
2. How many contigs are there?
3. How many features are on the positive strand?
4. How many features are on the negative strand?
5. How many repeat features are there?

When you have completed these exercises move on to the next part of the tutorial, [Advanced Linux](#).

4 Advanced Linux

So far, we have run single commands in a terminal. Now we will look at ways to run more than a single command and how to automate tasks.

Before you start this section change into the `advanced_linux` directory:

```
 cd ../advanced_linux
```

4.1 Repeating analysis with loops


It is common in Bioinformatics to run the same analysis on many files. Suppose we had a command that ran one type of analysis, and we wanted to repeat the same analysis on 100 different files. It would be tedious, and error-prone, to write the same command 100 times. Instead we can use a loop.

As an example, say we wanted to run `wc` on each file in the directory `loop_files`.

First let's look at the contents of the `loop_files` directory:

```
 ls loop_files/
```


To run `wc` on each of the files found in the directory `loop_files/` use the following command:

```
 for filename in loop_files/*; do wc $filename; done
```

However, it is useful to be able to run multiple commands that process some data and produce some output. These commands can be put into a file (i.e. a script), and run on input data. This has the advantage of saving time and reproducibility, so that the same analysis can be run on many input data sets.

4.2 Introduction to BASH scripting

It is traditional when learning a new language (in this case BASH), to write a script that says “Hello World!”. Open a terminal and make a new directory in your home called `scripts`, by typing

```
 cd ~  
mkdir ~/scripts
```

Next open a text editor, which you will use to write the script. What text editors are available will depend on your system. For example, `gedit` in Linux. Do not try to use a word processor, such as Word! If you don't already have a favorite, try `nano` by running the following command:

```
 nano scripts/hello.sh
```

Type this into the text editor:

```
echo Hello World!
```


and save this to a file called `hello.sh` in your new `scripts` directory. This script will print `Hello World!` to the screen when we run it. First, check that the script is saved in the correct place.

```
ls scripts/hello.sh
```

Now try to run the script. For now, we need to tell Linux that this is a BASH script and where it is (inside the `scripts` directory):

```
bash scripts/hello.sh
```

4.3 Setting up a scripts directory

It would be nice if our scripts could be run from anywhere in the filesystem, without having to tell Linux where the script is, or that it is a BASH script. This is how built-in commands work, like `cd` or `ls`.

To tell Linux that the script is a BASH script, edit the file and add this line as the first line of the script:

```
#!/usr/bin/env bash
```

and remember to save the script again. This special line at the start of the file tells Linux that the file is a bash script, so that it expects bash commands throughout the file. There is one more change to be made to the file to tell Linux that it is a program to be run (it is “executable”). This is done with the command `chmod`. Type this into the terminal to make the file executable:

```
chmod +x scripts/hello.sh
```

Now, the script can be run, but we must still tell Linux where the script is in the filesystem. In this case, it is in a directory called `scripts` in the current working directory, “`./scripts`”.

```
./scripts/hello.sh
```

We need to change our setup so that Linux can find the script without us having to explicitly say where it is. Whenever a command is typed into Linux, it has a list of directories that it searches through to look for the command. To see the list of directories type:

```
echo $PATH
```

It returns a list of directories, which are all the places Linux will look for a command. First, check what happens if we try to run the script without telling Linux where it is:

```
hello.sh
bash: hello.sh: command not found
```

Linux did not find it! The command to run to add the `scripts` directory to `$PATH` is:

```
export PATH=$PATH:~/scripts/
```

If you want this change to be permanent, ie so that Linux finds your scripts after you restart your computer or logout and login, add that line to the end of a file called `~/.bashrc`. If you are using a Mac, then the file should instead be `~/.bash_profile`. If the file does not already exist, then create it and put that line into it.

Now the script works, no matter where we are in the filesystem. Linux will check the scripts directory and find the file `hello.sh`. You can be *anywhere* in your filesystem, and simply running

```
hello.sh
```

will always work. Try it now.



```
hello.sh
```

In general, when making a new script, you can now copy and edit an existing script, or make a new one like this:

```
cd ~/scripts
touch my_new_script.sh
chmod +x my_new_script.sh
```


and then open `my_new_script.sh` in a text editor.

If you would like to learn more advanced bash scripting we have provided some further optional material in [BASH Scripting](#).

5 BASH scripting

This section is additional material and provided for anyone who would like to learn more advanced bash scripting.

Before you start this section change into the `bash_scripts` directory:


```
 cd ../bash_scripts
```


5.1 Getting options from the terminal and printing a help message

Usually, we would like a script to read in options from the user, such as the name of an input file. This would mean a script can be run like this:

```
my_script.sh input_file
```

Inside the script, the parameters provided by the user are given the names `$1`, `$2`, `$3` etc (do not confuse these with column names used by `awk`!). Here is a simple example that expects the user to provide a filename and a number. The script simply prints the filename to the screen, and then the first few lines of the file (the number of lines is determined by the number given by the user).

```
 cat options_example.sh
```

```
 options_example.sh test_file 2
```

The options have been used by the script, but the script itself is not very readable. It is better to use names instead of `$1` and `$2`. Here is an improved version of the script that does exactly the same as the previous script, but is more readable.

```
 cat options_example.2.sh
```

5.2 Checking options from the user

The previous scripts will have strange behaviour if the input is not as expected by the script. Many things could go wrong. For example:

- The wrong number of options are given by the user
- The input file does not exist.

Try running the script with different options and see what happens.

A convention with scripts is that it should output a help message if it is not run correctly. This shows anyone how the script should be run (including you!) without having to look at the code inside the script.

A basic check for this script would be to verify that two options were supplied, and if not then print a help message. The code looks like this:

```
if [ $# -ne 2 ]
then
    echo "usage: options_example.3.sh filename number_of_lines"
```

```
    echo
    echo "Prints the filename, and the given first number of lines of the file"
    exit
fi
```

You can copy this code into the start of any of your scripts, and easily modify it to work for that script. A little explanation:

- A special variable `$#` has been used, which is the number of options that were given by the user.
- The whole block of code has the form “`if [$# -ne 2] then fi`”. This only runs the code between the `then` and `fi`, if `$#` (the number of options) is not 2.
- The line `exit` simply makes the script end, so that no more code is run.



```
options_example.3.sh
```

Another check is that the input file really does exist. If it does not exist, then there is no point in trying to run any more code. This can be checked with another `if ... then ... fi` block of code:

```
if [ ! -f $filename ]

then
    echo "File '$filename' not found! Cannot continue"
    exit
fi
```

Putting this all together, the script now looks like this:



```
cat options_example.3.sh
```

Two new features have also been introduced in this file:

1. The second line is “`set -eu`”. Without this line, if any line produces an error, the script will carry on regardless to the end of the script. Using the `-e` option, an error anywhere in the file will result in the script stopping at the line that produced the error, instead of continuing. In general, it is best that the script stops at any error. The `-u` creates an error if you try to use a variable which doesn’t exist. This helps to stop typos doing bad things to your analysis.
2. There are several lines starting with a hash `#`. These lines are “comment lines” that are not run. They are used to document the code, containing explanations of what is happening. It is good practice to comment your scripts!

The above script provides a template for writing your own scripts. The general method is:

1. Tell Linux that this is a BASH script, and to stop at the first error.
2. Check if the user ran the script correctly. If not, output a message telling the user how to run the script.
3. Check the input looks OK (in this case, that the input file exists).
4. Process the input.

5.3 Using variables to store output from commands


It can be useful to run a command and put the results into a variable. Recall that we stored the input from the user in sensibly named variables:


```
filename=$1
```

The part after the equals sign could actually be any command that returns some output. For example, running this in Linux


```
wc -l filename | awk '{print $1}'
```

returns the number of lines. In case you are wondering why the command includes `| awk '{print $1}'`, check what happens with and without the pipe to `awk`:

```
 wc -l options_example.3.sh
```

```
 wc -l options_example.3.sh | awk '{print $1}'
```

With a small change, this can be stored in a variable and then used later.

```
 filename=options_example.3.sh  
line_count=$(wc -l $filename | awk '{print $1}')
```

```
echo There are $line_count lines in the file $filename
```

5.4 Exercises

1. Write a script that gets a filename from the user. If the file exists, it prints a nice human-readable message telling the user how many lines are in the file.
2. Use a loop to run the script from Exercise 1 on the files in the directory `loop_files/`.
3. Write a script that takes a GFF filename as input. Make the script produce a summary of various properties of the file. There is an example input file provided called `bash_scripts/exercise_3.gff`. Use your imagination! You could have a look back at the `awk` section of the course for inspiration. Here are some ideas you may wish to try:
 - Does the file exist?
 - How many records (ie lines) are in the file?
 - How many genes are in the file?
 - Is the file badly formatted in any way (eg wrong number of columns, do the coordinates look like numbers)?

6 Linux Quick Reference Guide

6.1 Looking at files and moving them around

```
pwd # Tell me which directory I'm in
ls # What else is in this directory
ls .. # What is in the directory above me
ls foo/bar/ # What is inside the bar directory which is inside the foo/ directory
ls -lah foo/ # Give the the details (-l) of all files and folders (-
a) using human
      # readable file sizes (-h)
cd ../.. # Move up two directories
cd ../foo/bar # Move up one directory and down into the foo/bar/ subdirectories
cp -r foo/ baz/ # Copy the foo/ directory into the baz/ directory
mv baz/foo .. # Move the foo directory into the parent directory
rm -r ../foo # remove the directory called foo/ from the parent directory
find foo/ -name "*.gff" # find all the files with a gff extension in the directory foo/
```

6.2 Looking in files

```
less bar.bed # scroll through bar.bed
grep chrom bar.bed | less -S # Only look at lines in bar.bed which have 'chrom' and
      # don't wrap lines (-S)
head -20 bar.bed # show me the first 20 lines of bar.bed
tail -20 bar.bed # show me the last 20 lines
cat bar.bed # show me all of the lines (bad for big files)
wc -l bar.bed # how many lines are there
sort -k 2 -n bar.bed # sort by the second column in numerical order
```

6.3 Grep

```
grep foo bar.bed # show me the lines in bar.bed with 'foo' in them
grep foo baz/* # show me all examples of foo in the files immediately within baz/
grep -r foo baz/ # show me all examples of foo in baz/ and every subdirectory within it
grep '^foo' bar.bed # show me all of the lines beginning with foo
grep 'foo$' bar.bed # show me all of the lines ending in foo
grep -i '^[acgt]$_' bar.bed # show me all of the lines which only have the characters
      # a,c,g and t (ignoring their case)
grep -v foo bar.bed # don't show me any files with foo in them
```

6.4 Awk

```
awk '{print $1}' bar.bed # just the first column
awk '$4 ~ /^foo/' bar.bed # just rows where the 4th column starts with foo
awk '$4 == "foo" {print $1}' bar.bed # the first column of rows where the 4th column is foo
awk -F"\t" '{print $NF}' bar.bed # ignore spaces and print the last column
awk -F"\t" '{print $(NF-1)}' bar.bed # print the penultimate column
```

6.5 Piping, redirection and more advanced queries

```
grep -hv '^#' bar/*.gff | sort
# grep => -h: don't print file names
#         -v: don't give me matching files
#         '^#': get rid of the header rows
#         'bar/*.gff': only look in the gff files in bar/
# sort => sort results returned from grep
```

6.6 A script

```
#!/usr/bin/env bash

set -e # stop running the script if there are errors
set -u # stop running the script if it uses an unknown variable
set -x # print every line before you run it (useful for debugging but annoying)

if [ $# -ne 2 ]
then
    echo "You must provide two files"
    exit 1 # exit the programme (and number > 0 reports that this is a failure)
fi

file_one=$1
file_two=$2

if [ ! -f $file_one ]
then
    echo "The first file couldn't be found"
    exit 2
fi

if [ ! -f $file_two ]
then
    echo "The second file couldn't be found"
    exit 2
fi

# Get the lines which aren't headers,
# take the first column and return the unique values
number_of_contigs_in_one=$(awk '$1 !~ /^#/ {print $1}' $file_one | sort -u | wc -l)
number_of_contigs_in_two=$(awk '/^[^#]/ {print $1}' $file_two | sort -u | wc -l)

if [ $number_of_contigs_in_one -gt $number_of_contigs_in_two ]
then
    echo "The first file had more unique contigs than the second"
    exit
```

```
elif [ $number_of_contigs_in_one -lt $number_of_contigs_in_two ]
then
    echo "The second file had more unique contigs"
    exit
else
    echo "The two files had the same number of contigs"
    exit
fi
```

6.7 Pro tips

- Use tab completion - it will save you time!
- Always have a quick look at files with `less` or `head` to double check their format
- Watch out for data in headers and that you don't accidentally `grep` some if you don't want them
- Watch out for spaces, especially if you're using `awk`; if in doubt, use `-F"\t"`
- Regular expressions are wierd, build them up slowly bit by bit
- If you did something smart but can't remember what it was, try typing `history` and it might have a record
- `man the_name_of_a_command` often gives you help
- Google is normally better at giving examples (prioritise stackoverflow.com results, they're normally good)

6.8 Which tool should I use?

You should probable use `grep` if:

- you're looking for files which contain some specific text (e.g. `grep -r foo bar/`: look in all the files in `bar/` for any with the word 'foo')

You should probably use `awk` if:

- your data has columns
- you need to do simple maths

You should use `find` if:

- you know something about a file (like it's name or creation date) but not where it is
- you want a list of all the files in a subdirectory and its subdirectories etc.

You should write a script if:

- your code doesn't fit on one line
- it's doing something you might want to do again in 3 months
- you want someone else to be able to do it without asking loads of questions
- you're doing something sensitive (e.g. deleting loads of files)
- you're doing something lots of times

You should probably use `less` or `head`:

- always, you should always use `less` or `head` to check intermediary steps in your analysis