

Programming for Bioinformatics
BIOL 8803 B
September 21st, 2015

This week will be about getting more experience with regex. The commands and concepts for today are:

`grep` – Get regular expression, search input for matching lines

`sed` – Stream editor, alter incoming data

`awk` – Get certain columns/rows from a file

Instructions for submission:

- Please download the files from T-square and place it in ~/class/ex5/ in your VM
- **This exercise has an associated quiz (quiz week 3) to it that can be found on the T-square**
- You will be required to submit a brief solution sheet like the one shown before in the lecture week 1 as well as complete the quiz.
- Again – the solution sheet will be brief i.e., **at most** 8 pages long if needed be.

Exercises

Regular expressions is one of the most intimidating topics in bioinformatics but if you can master it, you can master the art of file processing. This week we will learn efficient ways of file processing using advanced techniques. The learning objective of this week, as it has been for the class, is not memorizing every single piece of code. No one does that. It's about building strong conceptual foundations that you can use to solve problems in your research.

Regardless of which field you will go to, file manipulations techniques will always be useful. If there is only one thing that you can learn from this class, let it be this.

We'll be happy to go over any concept that you are still struggling with in Thursday's lab session.

General regex and programming

- 1.) Create regular expressions for lines that match these things
 - a.) Any number that is a multiple of 5
 - b.) Exactly 5 characters
 - c.) Any letter followed by a number
 - d.) The first 6 columns of a BED file
 - e.) Write a regular expression that matches the first 3 bases in a DNA sequence
 - f.) Write a regular expression that matches the last 3 bases in a DNA sequence
 - g.) Write a regular expression that matches two numbers followed by 2 lower case letters
 - h.) What does this regular expression match: `\d*\.\d{3}`

2.) What is the value of `a_number` after each line? (Write the value of `a_number` after each statement. These statements are in continuation)

```
a_number = 6
a_number -= 12
a_number += 15
a_number /= 3
a_number *= 2
a_number = a_number + a_number / 2
```

3.) What do these Boolean statements evaluate to? (Answer as simple as TRUE or FALSE)

- a.) `1 > 2`
- b.) `2 > 1`
- c.) `!(2 > 1)`
- d.) `(2 > 1) || (1 > 2)`
- e.) `(2 > 1) && (1 > 2)`
- f.) `!(2 > 1) || !(1 > 2)`
- g.) `!((2 > 1) || (1 > 2))`
- h.) `!((2 > 1) && (1 > 2))`

Regular Expression Exercises

6.) Searching a file with `grep`

- a.) Extract the `knownGene.txt.gz`. Google the command.
- b.) Use `grep` to get all genes on `chr22`
- c.) Use `grep` to get all and only those genes that occur on `chr1`

7.) Editing data streams with `sed`

- a.) Take the result of 6b and duplicate each line
- b.) Change the `'chr'` portion of every other line to `'cow'`
- c.) Delete the lines that have `cow` in them
- d.) Do a-c again, but this time do it 'in-place' in the file.

Real meat

4.) Format recognition

A lot of bioinformatics is all about data and the numerous formats that defines that data. Sooner or later you will have to be acquainted with these formats. Let's start today.

I understand there is way too many formats but honestly no one remembers all of them. You can always look them up.

We will deal with 4 file formats this week: GenBank, EMBL, FASTQ and MEGA. Here are their description:

EMBL	http://www.bioperl.org/wiki/EMBL_sequence_format
FASTQ	http://www.bioperl.org/wiki/FASTQ_sequence_format
GenBank	http://www.bioperl.org/wiki/GenBank_sequence_format
MEGA	http://www.bioperl.org/wiki/MEGA_multiple_alignment_format

We will learn more about this later in the course so don't worry a lot about them. Let's write simple `awk` one liners to check what format our input is in. This `awk` script will be provided only the first line off the file and it will predict what format the file is in.

You can send the first line using `head` command.

Format recognition can be made based on the first column itself.

The solution can be very clean and simple to something extremely complicated. It's all about how clearly you can think.

Files are provided on T-square for you to check your script.

5.) An *in silico* restriction enzyme digestion

In a parallel universe, restriction enzymes comes by the name of `sed` and breaks microbial genome using substitute patterns. One such enzyme has magically found its way to your computer. Download the M07149.fasta from T-square. We got some cutting to do.

a.) This restriction enzyme works on the pattern GAATTC and cuts it right after G like this:



```
GAATTC
CTTAAG
```

Cut the genome into pieces using this restriction enzyme (`sed`)! Store the fragmented genome in a new file. How many pieces did you get? (Don't count by hand. Use something sophisticated like `wc`)

b.) On further investigation, you found that the restriction enzyme is a little flexible. It can actually cut after the first base in the following patterns:

GAATTC, GAATTG,
GATTC, GATTTG,
CAATTC, CAATTG,
CATTC, CATTTG

Update your pattern to cut the genome accordingly. How many pieces did you get this time?

c.) You underestimated the strength of this enzyme. It can also vary its length. The updated list of pattern has the following letters optional: third (A or T), fourth (T) and last (C or G). Update the pattern to get the new number of pieces. How many do you get this time?

6.) Using `awk` to manipulate text on the command line

You should now be somewhat familiar with `awk`. Going over it again, `awk` is a very simple scripting language designed to work with text. It can be used to filter text, format text, gather information about the text, etc. I use it when I want to do something very simple. If I start getting into `awk` scripts that are longer than a few lines, I give up and do it in perl. Use `awk` to do the following with a BED file.

Download the **UCSC refGene.txt file** in standard BED format i.e., 12 columns. To do this, go to UCSC genome browser (<http://genome.ucsc.edu/>). You should see Table Browser in the left panel. Click on it. Select the following options:

clade: genome: assembly:

group: track:

table:

region: ☒ genome ☐ ENCODE Pilot regions ☐ position

identifiers (names/accessions):

filter:

intersection:

correlation:

output format: Send output to ☐ [Galaxy](#) ☐ [GREAT](#) ☐ [GenomeSpace](#)

output file: (leave blank to keep output in browser)

file type returned: ☒ plain text ☐ gzip compressed

Hold on to this original file. I am going to refer to this as **ref.bed**. You can call it whatever you like. I will keep asking you to come back to this. For questions where you have to edit this file, I will always ask you to save the edited version as a new file.

a.) Starting simple, print the unaltered contents of the BED file.

b.) Print **every other** line in the BED file.

c.) Print out entries that are on the negative strand only.

d.) For entries on the negative strand, print out the **coordinates** (not the file order) in reverse order – some programs want it like this for reasons I’ve never really understood. E.g.: if the input looked like this:

```
chr1 25071759 25170815 NM_013943 0 +
chr1 48998526 50489626 NM_032785 0 -
```

I want the output to be like this:

```
chr1 25071759 25170815 NM_013943 0 +
chr1 50489626 48998526 NM_032785 0 -
```

e.) Sum the total length of all of the entries in the BED file, i.e. from start to stop of each entry.

f.) Repeat previous question for those entries on the ‘+’ strand only.

g.) Change the ID (column 4) of each entry to its order of occurrence in the file, i.e. 0-n.

h.) Add a new column onto the end of the BED with the value of the new column being equal to the number of columns in the BED. Save this as a new file.

i.) Change the ‘chr’ portion of all entries on a single chromosome of your choice to ‘chicken’, leaving the others untouched.

7.) Let’s be awesome!

Download these files from t-square (and course repository <http://jordan.biology.gatech.edu/biol8803b/> if file is not available on t-square)

a.) Count the number of words in file1.txt using awk alone (no wc)

b.) Write an awk script to perform FASTQ to FASTA conversion. The FASTQ file is r2.fastq

c.) Write an awk script to perform SAM to FASTQ conversion. The SAM file is map.sam

d.) Write an awk script to extract reads with flag 12 and output it as a FASTQ file. The SAM file is map.sam

8.) Merging outputs

cat and paste aren’t the only commands capable of merging things. There is &&

a.) Let's start simple, merge the outputs of two `echo` commands using `&&`. This can be done as follows:
`echo "This is my first line" && echo "This is my second line"`

b.) Now let's get fancy. Print me the first 5 followed by the last 5 lines of `ref.bed`.

c.) Command inception: Print me lines the 5 lines that comes after the first 5 lines (i.e., only lines 6-10) followed by the 5 lines that comes right before the last 5 lines (i.e., say there are 100 lines, I need lines 90-95). Hint: you will require `()`.