

Command Line for Data Filtering Answers

Before you begin, copy “Oxy.counts.tsv” from “penelopeprime” to your home directory and then open Terminal. Alternatively, you can download the example file from:

<https://funcgen2022.buschlab.org/downloads/Oxy.counts.tsv>

1. Using the `awk` and `wc` commands (and a pipe), find out how many genes are significantly differentially expressed (i.e. adjusted p-value < 0.05).

```
awk '$3 < 0.05' Oxy.counts.tsv | wc -l  
341
```

2. Using the `cut` command, make a new file that just contains the Ensembl ID, the adjusted p-value, the \log_2 fold change and the gene name and description.

```
cut -f1,3,4,10,11 Oxy.counts.tsv > q2.tsv
```

3. Search for all the genes whose name begins with “si:”. How many are there?

```
cut -f10 Oxy.counts.tsv | grep si: | wc -l  
3558
```

4. How many genes have a biotype of “protein_coding”?

```
cut -f9 Oxy.counts.tsv | grep protein_coding | wc  
-l  
25432
```

5. Using just the `awk` command, make a new file that contains the Ensembl ID, gene name, chromosome and strand (in that order) for all the genes on the reverse strand.

```
awk -F"\t" '$8 == -1 { print $1 "\t" $10 "\t" $5  
"\t" $8 }' Oxy.counts.tsv > q5.txt
```

6. Use the `man` command to find out about the `more` command. What option do you need to use with `more` to see line numbers in the `Oxy.counts.tsv` file? (If the `more` command doesn't have a suitable option then look at the `less` command instead. `less` is an updated version of `more`.)

It turns out the version of `more` on the training room computers doesn't have an option for this. Instead, you need to use `less`.

```
less -N Oxy.counts.txt
```

7. Use the `sort` command to order the file by chromosome. Does the order of the non-numeric "chromosomes" make sense? Try using the `-V` option of `sort`, instead of `-g`. Is the order now better? (The `-V` option is technically for sorting version numbers, but it's also really useful for sorting chromosome names!)

```
sort -g -k5 Oxy.counts.tsv | more
```

```
sort -V -k5 Oxy.counts.tsv | more
```

8. How many genes are between 10,000,000 bp and 20,000,000 bp on chromosome 1?

```
awk '$5 == "2" && $7 > 10000000 && $6 < 20000000'  
Oxy.counts.tsv | wc -l  
197
```