# BIO720_Week2_hr1

Ian Dworkin

2024-09-10

## Contents

# In class activities for week 2

Now that you have had a bit of an introduction to the fundamentals of R, let's start working with some examples of how we can use these skills with real data.

We are going to work with the count data of RNAseq sequence reads from a large data set from a paper we recently submitted. It was part of a study of the genomic basis of the evolutionary of sociability in *Drosophila melanogaster".

The link to the data and script repository is here

But for the moment, all we need is the counts that were generated. For this particular analysis this was done by aligning the sequence reads to the *Drosophila melanogaster* genome using a *splice aware* alignment tool, STAR. While it is not relevant for the moment, if you want to read about this approach, you can take a look at this tutorial.

In any case, the count data can be downloaded like this, directly from our github repository.

```
rna_counts <- read.table("https://raw.githubusercontent.com/DworkinLab/DrosophilaSociabilityTranscript
                         header = TRUE)
```

## What we did yesterday (Sept 10th, 2024)

We first thought through some things we should check after we downloaded the data

**what should we look at first?**

- does the data look like it is expected
  - size of the data. Dimensions of the data?
  - names of variables?

- how is missing data encoded?
- what object class are we working with? Is it what we expected?
- Is all the data there? how do we check
- are column headers being read as headers?
  * Is the data possible?

We realized there were a variety of ways to check the number of rows (and columns of the data

Probably the most concise

```
dim(rna_counts)
```

Has lots of other useful information, but maybe too much information in this case!

```
str(rna_counts)
```

We can ask about how many rows and how many columns, *per se*

```
nrow(rna_counts)
ncol(rna_counts)
```

We then checked to see if there was any missing data (encoded by `NA`)

```
anyNA(rna_counts)
```

**let's finish answering the questions we had regarding what kind of object we have**

## Let's look at the names of the variables we have

So what are we looking at?

**A small but important change**

Typically in R and bioconductor for genomics experiments, we have columns as samples and rows as features (transcripts, genes, SNPs, taxa, etc). How is it currently set up?

How would we make a new object (generally don't over-write the old one!) where it corresponds to this format? Write down the steps you need to take in order to do so (the pseudo-code)

**extract the information about the experimental design from the column names for each sample**

We used underscored "_" for our delimiter in the files

- AS: Initials of who collected the samples
- the next is for the evolutionary treatments
- Replicate lineage
- Sex
- Environment
- Unique sample
- Lane of Illumina flowcell

## Pseudo code

Write pseudo-code to take the sample names and generate an object (but what kind of object?) to store all of the experimental meta-data.