

Introduction to Bulk RNAseq data analysis

Differential Expression of RNA-seq data

Last modified: 30 Sep 2024

Contents

| | |
|---|-----------|
| 1. R packages for differential gene expression analysis | 1 |
| 2. Load the counts data and sample meta data | 2 |
| 3. Creating the design model formula | 2 |
| The simple model | 2 |
| Exercise 1 | 3 |
| 4. Build a DESeq2DataSet | 3 |
| Filter out the unexpressed genes | 4 |
| 5. Differential expression analysis with DESeq2 | 4 |
| The DESeq2 work flow | 4 |
| The DESeq command | 7 |
| 6. Generate a table of differential expression results | 7 |
| Exercise 2 | 8 |
| Independent filtering | 8 |
| 7. The additive model | 9 |
| Exercise 3 | 9 |
| Exercise 4 | 10 |
| Exercise 5 | 10 |
| 8. The interaction model | 10 |
| Exercise 6 | 12 |
| 9. Extracting specific contrasts from an interaction model | 12 |
| Exercise 7 | 13 |
| 10. Save the results | 13 |

1. R packages for differential gene expression analysis

Now that we are happy that the quality of the data looks good, we can proceed to testing for differentially expressed genes. There are a number of packages to analyse RNA-Seq data. Most people use DESeq2 (Love, Huber, and Anders 2014) or edgeR (Robinson, McCarthy, and Smyth 2010; McCarthy, Chen, and Smyth 2012). There is also the option to use the limma package and transform the counts using its `voom` function. They are all equally valid approaches (Ritchie et al. 2015). There is an informative and honest blog post here by Mike Love, one of the authors of DESeq2, about deciding which to use.

We will use **DESeq2** for the rest of this practical.

To start with we will need three pieces of information to pass to DESeq2:

- The raw counts data
- The sample meta data
- The design model

2. Load the counts data and sample meta data

In the previous session we read the results from Salmon into R and created a `txi` object, which we then saved into an “rds” file. We can now load the `txi` from that file to start the differential expression analysis. We will also need the sample meta data sheet

First load the packages we need.

```
library(DESeq2)
library(tidyverse)
```

Now load the data from the earlier session.

```
txi <- readRDS("RObjects/txi.rds")
sampleinfo <- read_tsv("data/samplesheet_corrected.tsv", col_types = "cccc")
```

It is important to be sure that the order of the samples in rows in the sample meta data table matches the order of the columns in the data matrix - `DESeq2` will **not** check this. If the order does not match you will not be running the analyses that you think you are.

```
all(colnames(txr$counts) == sampleinfo$SampleName)
## [1] TRUE
```

3. Creating the design model formula

Next we need to create a design model formula for our analysis. `DESeq2` will use this to generate the model matrix, as we have seen in the linear models lecture.

We have two variables in our experiment: “Status” and “Time Point”.

We will fit two models under two assumptions: no interaction and interaction of these two factors, however, to demonstrate the how `DESeq2` is used we will start with a simple model which considers just one factor.

The simple model

First, create a variable containing the model using standard R ‘formula’ syntax.

```
simple.model <- as.formula(~ TimePoint)
```

What does this look like as a model matrix?

```
model.matrix(simple.model, data = sampleinfo)
```

```
## (Intercept) TimePointd33
## 1           1          0
## 2           1          0
## 3           1          0
## 4           1          1
## 5           1          1
## 6           1          1
## 7           1          0
## 8           1          0
```

```

## 9      1      0
## 10     1      1
## 11     1      1
## 12     1      1
## attr(),"assign")
## [1] 0 1
## attr(),"contrasts")
## attr(),"contrasts")$TimePoint
## [1] "contr.treatment"

```

The model matrix haStatuss two columns. We can see that our variable “TimePoint” has been turned into two indicator variables called `(Intercept)` and `TimePointd33`. These match what we talked about earlier:

- `(Intercept)` is the indicator variable for β_0 . This is our reference group, which in this case is day 11.
- `TimePointd33` is the indicator variable for β_1 , and tells us whether samples belong to the day 33 group.

Thus, our β_1 coefficient is the difference in gene expression between day 33 and day 11.

Exercise 1

This time create and investigate the model matrix for the variable “Status”. 1. Create a model formula to investigate the effect of “Status” on gene expression. 2. Look at the model matrix and identify which is the reference group in your model.

The intercept has been set automatically to the group in the factor that is alphabetically first: `Infected`.

It would be nice if `Uninfected` were the base line/intercept. To get R to use `Uninfected` as the intercept we need to use a `factor`. Let’s set factor levels on Status to use `Uninfected` as the intercept.

```

sampleinfo <- mutate(sampleinfo, Status = fct_relevel(Status, "Uninfected"))
model.matrix(simple.model, data = sampleinfo)

```

```

##   (Intercept) TimePointd33
## 1      1      0
## 2      1      0
## 3      1      0
## 4      1      1
## 5      1      1
## 6      1      1
## 7      1      0
## 8      1      0
## 9      1      0
## 10     1      1
## 11     1      1
## 12     1      1
## attr(),"assign")
## [1] 0 1
## attr(),"contrasts")
## attr(),"contrasts")$TimePoint
## [1] "contr.treatment"

```

4. Build a DESeq2DataSet

We don’t actually need to pass `DESeq2` the model matrix, instead we pass it the design formula and the `sampleinfo` it will build the matrix itself. We will continue to use the simple model for `Status` that we just created. If you haven’t already done that, do so now:

```
simple.model <- as.formula(~ Status)
sampleinfo <- mutate(sampleinfo, Status = fct_relevel(Status, "Uninfected"))
```

We can now build our DESeq2 object using the three necessary components:

- The `txi` object containing the counts
- The `sampleinfo` data frame containing the sample metadata
- The `simple.model` design formula

```
# create the DESeqDataSet object
ddsObj.raw <- DESeqDataSetFromTximport(tx = txi,
                                         colData = sampleinfo,
                                         design = simple.model)
```

```
## using counts and average transcript lengths from tximport
```

When we summarised the counts to gene level, `tximport` also calculated an average transcript length for each gene for each sample. For a given gene the average transcript length may vary between samples if different samples are using alternative transcripts. DESeq2 will incorporate this into its “normalisation”.

Filter out the unexpressed genes

Just as we did in session 7, we should filter out genes that uninformative.

```
keep <- rowSums(counts(ddsObj.raw)) > 5
ddsObjfilt <- ddsObj.raw[keep,]
```

5. Differential expression analysis with DESeq2

The DESeq2 work flow

The main DESeq2 work flow is carried out in 3 steps:

`estimateSizeFactors`

First, Calculate the “median ratio” normalisation size factors for each sample and adjust for average transcript length on a per gene per sample basis.

```
ddsObj <- estimateSizeFactors(ddsObj.filt)
```

```
## using 'avgTxLength' from assays(dds), correcting for library size
```

Let's have a look at what that did DESeq2 has calculated a normalizsation factor for each gene for each sample.

```
normalizationFactors(ddsObj.filt)
```

```
## NULL
```

```
normalizationFactors(ddsObj)
```

```
##          SRR7657878 SRR7657881 SRR7657880 SRR7657874 SRR7657882
## ENSMUSG000000000001 0.9650391 0.96894491 0.94927826 0.9018189 1.20992201
## ENSMUSG000000000028 1.1565743 1.00859533 0.94221210 0.8798304 1.22726399
## ENSMUSG000000000037 0.9372735 1.11991330 0.52624731 1.0355843 0.77648936
##          SRR7657872 SRR7657877 SRR7657876 SRR7657879 SRR7657883
## ENSMUSG000000000001 0.9615574 1.1078271 1.01950314 0.95990882 0.9028895
## ENSMUSG000000000028 1.0086842 1.2535886 0.95400723 0.77315414 0.7550835
```

```

## ENSMUSG000000000037 0.8718128 2.2493164 0.77345594 0.67601042 0.5549084
##                               SRR7657873 SRR7657875
## ENSMUSG000000000001 1.0472683 1.04743598
## ENSMUSG000000000028 1.0609078 1.12777706
## ENSMUSG000000000037 1.8562184 2.13169578
## [ reached getOption("max.print") -- omitted 20088 rows ]

```

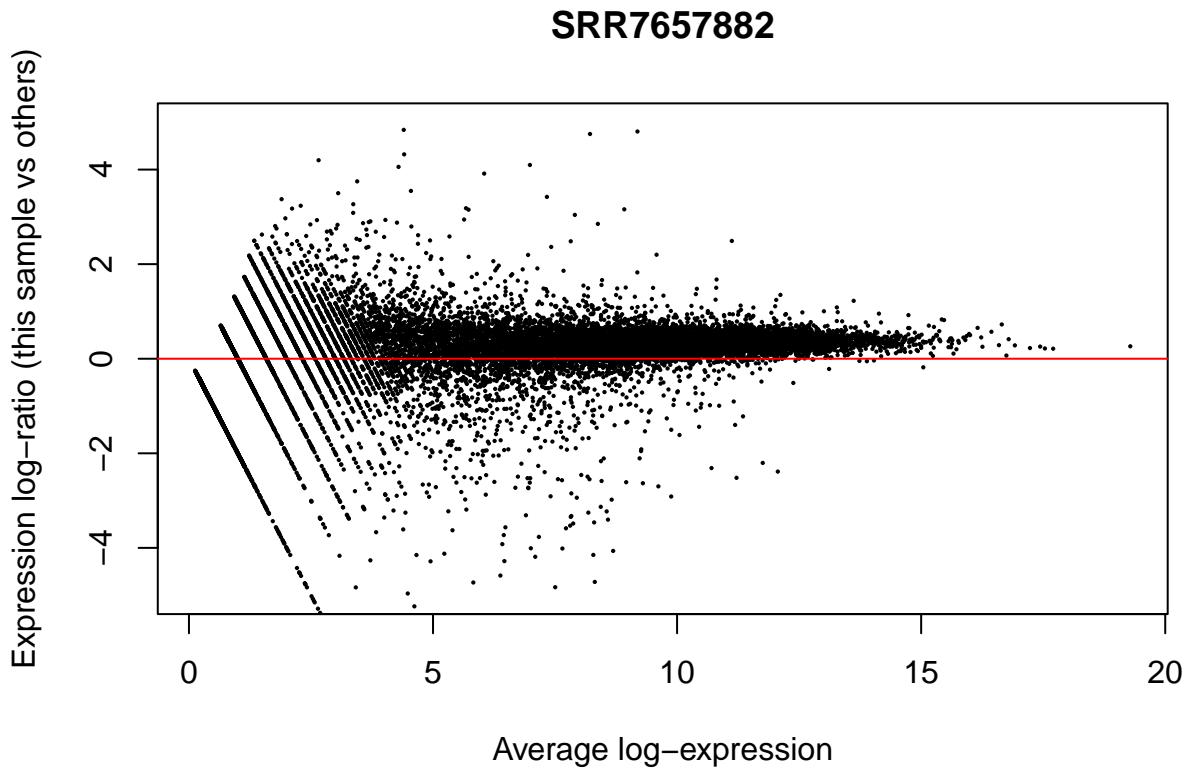
We can use `plotMA` from `limma` to look at the effect of these normalisation factors on data in an MA plot. Let's look at **SRR7657882**, the fifth column, which has the largest normalisation factors.

```

logcounts <- log2(counts(dds0bj, normalized = FALSE) + 1)

limma::plotMA(logcounts, array = 5, ylim = c(-5, 5))
abline(h = 0, col = "red")

```



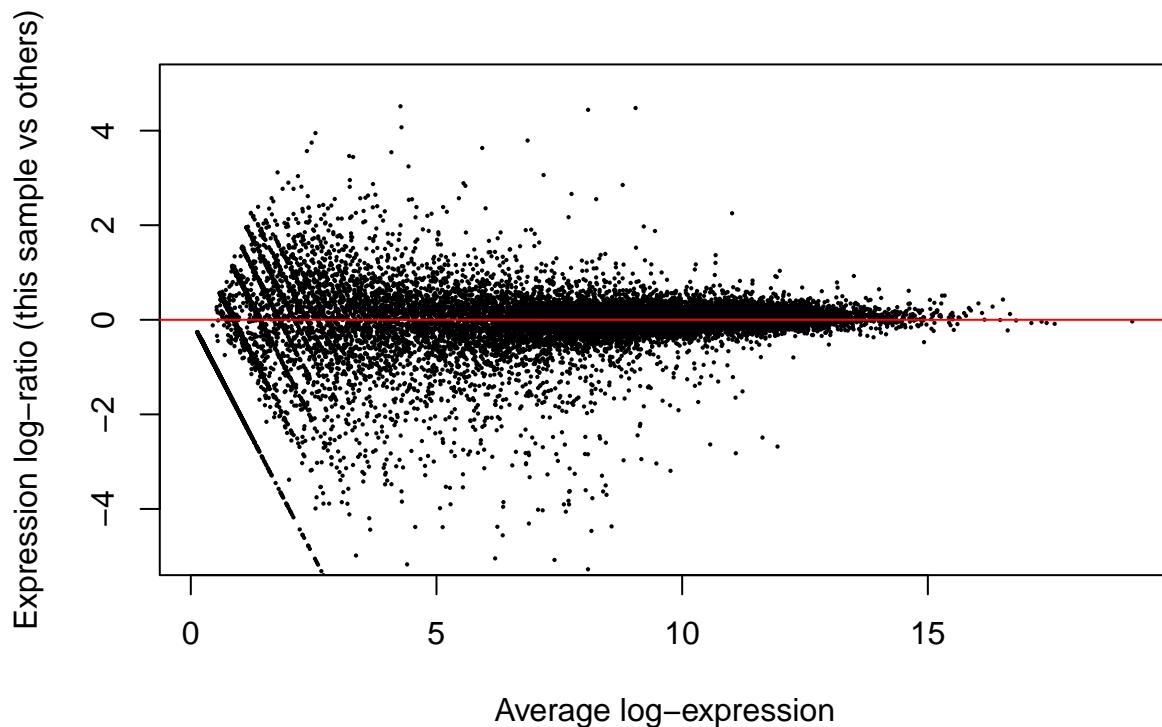
```

logNormalizedCounts <- log2(counts(dds0bj, normalized = TRUE) + 1)

limma::plotMA(logNormalizedCounts, array = 5, ylim = c(-5, 5))
abline(h = 0, col = "red")

```

SRR7657882



DESeq2 doesn't actually normalise the counts, it uses raw counts and includes the normalisation factors in the modeling as an “offset”. Please see the DESeq2 documentation if you'd like more details on exactly how they are incorporated into the algorithm. For practical purposes we can think of it as a normalisation.

`estimateDispersions`

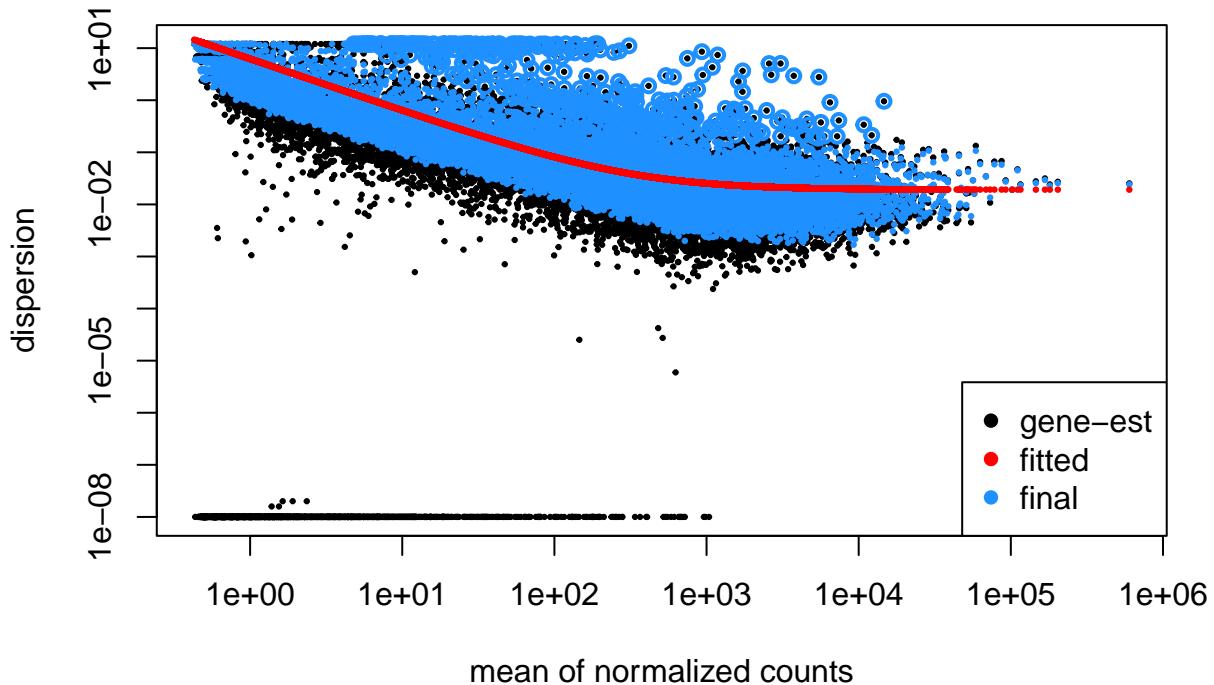
Next we need to estimate the dispersion parameters for each gene.

```
ddSObj <- estimateDispersions(ddSObj)
```

```
## gene-wise dispersion estimates  
## mean-dispersion relationship  
## final dispersion estimates
```

We can plot all three sets of dispersion estimates. It is particularly important to do this if you change any of the default parameters for this step.

```
plotDispEsts(ddSObj)
```



nbinomWaldTest

Finally, apply Negative Binomial GLM fitting and calculate Wald statistics.

```
ddsObj <- nbinomWaldTest(ddsObj)
```

The DESeq command

In practice the 3 steps above can be performed in a single step using the `DESeq` wrapper function. Performing the three steps separately is useful if you wish to alter the default parameters of one or more steps, otherwise the `DESeq` function is fine.

```
ddsObj <- DESeq(ddsObj.filt)
```

```
## estimating size factors
## using 'avgTxLength' from assays(dds), correcting for library size
## estimating dispersions
## gene-wise dispersion estimates
## mean-dispersion relationship
## final dispersion estimates
## fitting model and testing
```

6. Generate a table of differential expression results

We can generate a table of differential expression results from the DDS object using the `results` function of `DESeq2`.

```
results.simple <- results(ddsObj, alpha = 0.05)
results.simple
```

```

## log2 fold change (MLE): Status Infected vs Uninfected
## Wald test p-value: Status Infected vs Uninfected
## DataFrame with 20091 rows and 6 columns
##           baseMean log2FoldChange      lfcSE       stat      pvalue
##           <numeric>     <numeric> <numeric>     <numeric>     <numeric>
## ENSMUSG000000000001 1102.56094    -0.00802952  0.102877   -0.078050   0.937788
## ENSMUSG000000000028  58.60055     0.30498077  0.254312    1.199239   0.230435
## ENSMUSG000000000037  49.23586     -0.05272685  0.416862   -0.126485   0.899348
## ENSMUSG000000000049  7.98789     0.38165132  0.644869    0.591827   0.553966
## ENSMUSG000000000056 1981.00402    -0.16921845  0.128542   -1.316449   0.188024
##          padj
##          <numeric>
## ENSMUSG000000000001  0.975584
## ENSMUSG000000000028  0.480598
## ENSMUSG000000000037  0.961314
## ENSMUSG000000000049  0.772123
## ENSMUSG000000000056  0.426492
## [ reached getOption("max.print") -- omitted 6 rows ]

```

Exercise 2

Now we have made our results table using our simple model, let have a look at which genes are changing and how many pass our 0.05 threshold.

- how many genes are significantly (with an FDR < 0.05) up-regulated?
- how many genes are significantly (with an FDR < 0.05) down-regulated?
- Here is the results table for two of the genes:

| | baseMean | log2FoldChange | lfcSE | stat | pvalue | padj |
|--------------------|----------|----------------|-------|-------|---------|--------|
| ENSMUSG00000053747 | 34.930 | -1.61 | 0.531 | -3.04 | 0.00238 | 0.0199 |
| ENSMUSG00000048763 | 30.853 | -2.46 | 1.450 | -1.70 | 0.08970 | 0.2700 |

One of these is strongly downregulated with $\log_2(\text{fold-change})$ of -2.46 . On a linear scale this is $2^{-2.46} = 5.5$ times more highly expressed in the uninfected group relative to the infected group, but its adjusted p-value is 0.27. By contrast, the other gene has a lower LFC ~ -1.61 (on a linear scale $2^{-1.61} = 3.05$), but it's adjusted p-value is 0.0023.

How can you explain this apparent contradiction?

Independent filtering

From `DESeq2` manual: “The results function of the `DESeq2` package performs independent filtering by default using the mean of normalized counts as a filter statistic. A threshold on the filter statistic is found which optimizes the number of adjusted p values lower than a [specified] significance level”.

The default significance level for independent filtering is 0.1, however, you should set this to the FDR cut off you are planning to use. We will use 0.05 - this was the purpose of the `alpha` argument in the previous command.

Remember in Session 7 we said that there is no need to pre-filter the genes as `DESeq2` will do this through a process it calls ‘independent filtering’. The genes with NA are the ones `DESeq2` has filtered out.

7. The additive model

So far we have fitted a simple model considering just “Status”, but in reality we want to model the effects of both “Status” and “Time Point”.

Let’s start with the model with only main effects - an additive model with no interaction. The main assumption here is that the effects of Status and the effects of Time Point are independent.

First we create the additive model as a formula object.

```
additive.model <- as.formula(~ TimePoint + Status)
```

As before, we can now build the DESeq from the raw data, the sample meta data and the model, and then filter the data set.

```
ddsObj.raw <- DESeqDataSetFromTximport(tx = txi,
                                         colData = sampleinfo,
                                         design = additive.model)
keep <- rowSums(counts(ddsObj.raw)) > 5
ddsObj.filt <- ddsObj.raw[keep, ]
```

Exercise 3

You are now ready to run the differential gene expression analysis Run the DESeq2 analysis

1. Run the size factor estimation, dispersion estimation and modelling steps using the `DESeq` command as above.
2. Extract the default contrast using the `results` command into a new object called `results.additive`

Questions: a) How many coefficients are there in the additive model? b) What is the reference group in the additive model? c) What contrasts could we perform using this model? d) What contrast does the `results.additive` object represent? e) How many genes have an adjusted p-value of less than 0.05

The default contrast of `results` function

The `results` function has returned the results for the contrast “Infected vs Uninfected”. Let’s have a look at the model matrix to understand why `DESeq2` has given us this particular contrast.

```
model.matrix(additive.model, data = sampleinfo)
```

```
##   (Intercept) TimePointd33 StatusInfected
## 1           1          0           1
## 2           1          0           1
## 3           1          0           1
## 4           1          1           1
## 5           1          1           0
## 6           1          1           1
## 7           1          0           0
## 8           1          0           0
## 9           1          0           0
## 10          1          1           0
## 11          1          1           1
## 12          1          1           0
## attr(,"assign")
## [1] 0 1 2
## attr(,"contrasts")
```

```

## attr(,"contrasts")$TimePoint
## [1] "contr.treatment"
##
## attr(,"contrasts")$Status
## [1] "contr.treatment"

```

By default, `results` has returned the contrast encoded by the final column in the model matrix. `DESeq2` has the command `resultsNames` that allows us to view the contrasts that are available directly from the coefficients of the model.

```

resultsNames(ddsObj)

## [1] "Intercept"                  "TimePoint_d33_vs_d11"
## [3] "Status_Infected_vs_Uninfected"

```

Exercise 4

How do the named coefficients above relate to β_i coefficients in the design formula:

$$expression = \beta_0 + \beta_1 TimePoint_{d33} + \beta_2 Status_{Infected}$$

Let's just rename `results.additive` so that we know which contrast results it contains.

```

results.InfectedvUninfected <- results.additive
rm(results.additive)

```

Let's get the top 100 genes by adjusted p-value

```

topGenesIvU <- as.data.frame(results.InfectedvUninfected) %>%
  rownames_to_column("GeneID") %>%
  top_n(100, wt = -padj)
topGenesIvU

```

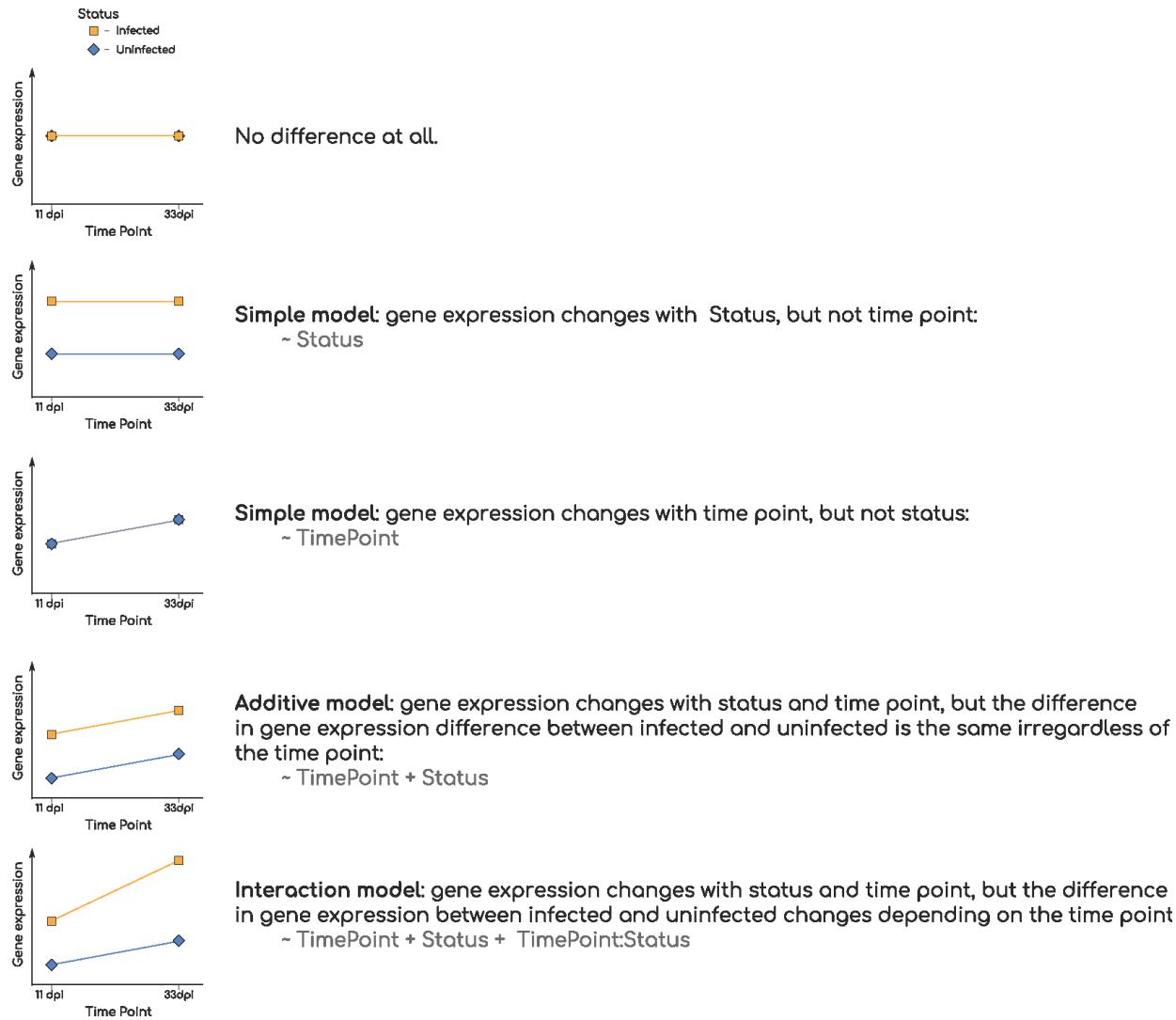
Exercise 5

If we want a different contrast we can just pass the `results` function the `name` of the contrast, as given by `resultsNames(ddsObj)`. Look at the help page for the `results` command to see how to do this.

1. Retrieve the results for the contrast of d33 versus d11.
2. How many differentially expressed genes are there at FDR < 0.05?

8. The interaction model

So far we have modeled gene expression as a function of Status and Time Point with an additive model. Realistically, we would probably expect the two factors interact such that differences in gene expression between infected and uninfected mice are not the same at different time points.



pdf version

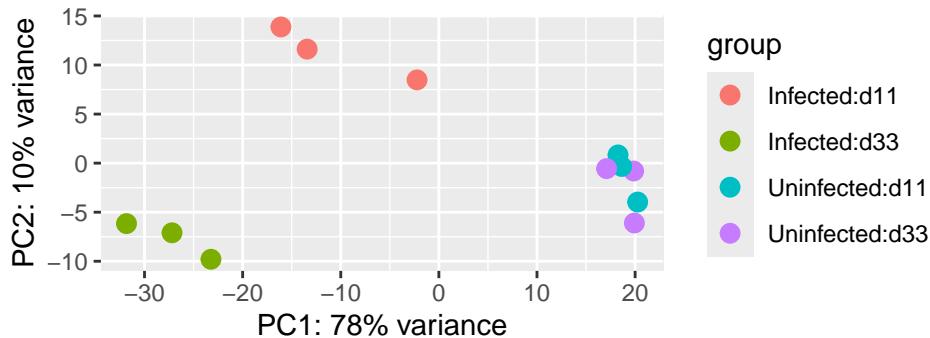
Let's plot a PCA from vst transformed data. Can you anticipate if the interaction term will be important?

```
vstcounts <- vst(dds0obj.raw, blind = TRUE)
```

```
## using 'avgTxLength' from assays(dds), correcting for library size
```

```
plotPCA(vstcounts, intgroup = c("Status", "TimePoint"))
```

```
## using ntop=500 top features by variance
```



In this case we can, from both the PCA and our understanding of the biology, be fairly certain that the interaction model is the appropriate model to use.

Exercise 6

1. Create a new DESeq2 object using a model with an interaction between TimePoint and Status. The model formula should be

```
~TimePoint + Status + TimePoint:Status
```

where `TimePoint:Status` is the parameter for the interaction beteween TimePoint and Status.

Note that `*` can be used as shortcut to add the interaction term, e.g. `~TimePoint * Status`, however, writing out in long form is clearer here.

Remember to filter to remove uninformative genes.

2. Run the statistical analysis using the `DESeq` command and create a new analysis object called `ddsObj.interaction`.
3. Extract a table of results using the default `results` command. What is the contrast that these results are for?

9. Extracting specific contrasts from an interaction model

With an interaction model we can no longer ask the general question “What is the difference in gene expression between Infected and Uninfected?”, but must rather ask two more specific questions:

- “What is the difference in gene expression between Infected and Uninfected at 11 days post infection?”
- “What is the difference in gene expression between Infected and Uninfected at 33 days post infection?”

If we view the `resultsNames` for the interaction model, we can see the intercept is Uninfected and 11 days post infection:

```
resultsNames(ddsObj.interaction)
```

```
## [1] "Intercept"                  "TimePoint_d33_vs_d11"
## [3] "Status_Infected_vs_Uninfected" "TimePointd33.StatusInfected"
```

The main effect `Status_Infected_vs_Uninfected` is therefore the difference between Infected and Uninfected at 11 days post infection.

```
results.interaction.11 <- results(ddsObj.interaction,
                                    name = "Status_Infected_vs_Uninfected",
                                    alpha = 0.05)
```

To get the results for Infected versus Uninfected at 33 days post infection, we would need to add the interaction term `TimePointd33.StatusInfected`.

In the help page for `results` it shows us how to do this with a `contrast` in example 3.

```
results.interaction.33 <- results(ddsObj.interaction,
                                    contrast = list(c("Status_Infected_vs_Uninfected",
                                                      "TimePointd33.StatusInfected")),
                                    alpha = 0.05)
```

Number of genes with `padj < 0.05` for Test v Control at day 11:

```
sum(results.interaction.11$padj < 0.05, na.rm = TRUE)

## [1] 1072
```

Number of genes with `padj < 0.05` for Test v Control at day 33:

```
sum(results.interaction.33$padj < 0.05, na.rm = TRUE)

## [1] 2782
```

We can see that there is a strong difference in the effects of infection on gene expression between days 11 and 33.

Exercise 7

Let's investigate the uninfected mice

1. Extract the results for d33 v d11 for Infected mice.
How many genes have an adjusted p-value less than 0.05?
2. Extract the results for d33 v d11 for Uninfected mice.
How many genes have an adjusted p-value less than 0.05?
Is this remarkable?

Do these results suggest another approach to analysing this data set?

10. Save the results

Finally save the corrected sample metadata, the DESeq2 dataset, and the two DESeq2 results objects.

```
write_tsv(sampleinfo, "results/samplesheet_corrected.tsv")
saveRDS(ddsObj.interaction, "results/DESeqDataSet.interaction.rds")
saveRDS(results.interaction.11, "results/DESeqResults.interaction_d11.rds")
saveRDS(results.interaction.33, "results/DESeqResults.interaction_d33.rds")
```

Love, Michael I., Wolfgang Huber, and Simon Anders. 2014. “Moderated Estimation of Fold Change and Dispersion for RNA-Seq Data with DESeq2.” *Genome Biology* 15: 550. <https://doi.org/10.1186/s13059-014-0550-8>.

McCarthy, Davis J, Yunshun Chen, and Gordon K Smyth. 2012. “Differential Expression Analysis of Multifactor RNA-Seq Experiments with Respect to Biological Variation.” *Nucleic Acids Research* 40 (10): 4288–97. <https://doi.org/10.1093/nar/gks042>.

Ritchie, Matthew E, Belinda Phipson, Di Wu, Yifang Hu, Charity W Law, Wei Shi, and Gordon K Smyth. 2015. “limma Powers Differential Expression Analyses for RNA-Sequencing and Microarray Studies.” *Nucleic Acids Research* 43 (7): e47. <https://doi.org/10.1093/nar/gkv007>.

Robinson, Mark D, Davis J McCarthy, and Gordon K Smyth. 2010. “edgeR: A Bioconductor Package for Differential Expression Analysis of Digital Gene Expression Data.” *Bioinformatics* 26 (1): 139–40. <https://doi.org/10.1093/bioinformatics/btp616>.