

Identifying changes in RNA binding using a linear model

In this notebook, we will apply the linear model approach presented in (1_simple_example) to a real-life data set.

```
# load packages
library(tidyverse)
library(biobroom)
library(MSnbase)

# set up standardised plotting scheme
theme_set(theme_bw(base_size = 20) +
  theme(panel.grid.major=element_blank(),
    panel.grid.minor=element_blank(),
    aspect.ratio=1))

cbPalette <- c("#E69F00", "#56B4E9", "#009E73", "#F0E442", "#0072B2", "#D55E00", "#CC79A7", "#999999")
```

We start by reading in the data. Our input here is the protein-level quantification for the Nocodazole arrest/release experiment conducted for the OOPS NBT paper. In this experiment, we wanted to assess changes in RNA binding in arrested/released cells. To do this, we quantified “total” protein abundance and RNA-bound (extracted by OOPS) protein abundance. The peptide-level abundances have been aggregated to protein level abundance and center-median normalised. Proteins with missing values have been removed.

```
total_protein_quant <- readRDS("../raw/total_res_pro_agg_norm.rds")
rbp_protein_quant <- readRDS("../raw/rbp_res_pro_agg_norm.rds")
```

The input data are in MSnSets. As a reminder, the `MSnSet` class mimics the `ExpressionSet` class and contains 3 matrices: 1. assay data (obtained via: `exprs`) 2. feature data (`fData`) 3. phenotype data (`pData`)

The assay data is the quantification of the features (PSMs/peptides/proteins) and contains one column per sample

The feature data describes each feature, e.g peptide sequence, master protein accession, retention time etc etc

The phenotype data describes the samples

Let’s take a look at our total protein quantification data. If we “print” the object, we get a summary including the processing steps performed.

Here we have 2761 features (proteins) quantified across 10 samples.

The `varLabels` describe the “Condition”, “Replicate” and “Type” for these samples. We’ll take a look at these in more detail shortly.

We can see that there were original 20171 features which were combined into 18111 features using a user-defined function. This was the step at which peptides with the same sequence but different variable modifications were aggregated. Then, these 18111 features were combined into the 2761 features (peptide->protein aggregation). Finally, the data was center-median normalised and missing values with imputed with knn.

```
print(total_protein_quant)
```

```
## MSnSet (storageMode: lockedEnvironment)
```

```
## assayData: 2761 features, 10 samples
##   element names: exprs
## protocolData: none
## phenoData
##   sampleNames: Abundance.F1.126.Sample Abundance.F1.127N.Sample
##     ... Abundance.F1.131.Sample (10 total)
##   varLabels: Sample_name Condition Replicate Type
##   varMetadata: labelDescription
## featureData
##   featureNames: A0AVT1 AOMZ66 ... Q9Y6Y8 (2761 total)
##   fvarLabels: Checked Confidence ... CV.Abandance.F1.131.Sample
##     (47 total)
##   fvarMetadata: labelDescription
## experimentData: use 'experimentData(object)'
## Annotation:
## - - - Processing information - - -
## Subset [20171,10][20171,10] Thu Aug  9 16:17:13 2018
## Combined 20171 features into 18111 using user-defined function: Thu Aug  9 16:17:32 2018
## Combined 18111 features into 2761 using user-defined function: Thu Aug  9 16:17:46 2018
## Normalised (center.median): Thu Aug  9 16:17:51 2018
## Data imputation using knn Thu Aug  9 16:17:58 2018
##   Using default parameters
## MSnbase version: 2.4.2
```

Here's the top of the assay data

```
print(dim(exprs(total_protein_quant)))
```

```
## [1] 2761  10
```

```
print(head(exprs(total_protein_quant), 2))
```

```
##           Abundance.F1.126.Sample Abundance.F1.127N.Sample
## A0AVT1           5.387263           5.379489
## AOMZ66           5.018771           4.926415
##           Abundance.F1.127C.Sample Abundance.F1.128N.Sample
## A0AVT1           5.296457           5.356980
## AOMZ66           4.876517           4.955341
##           Abundance.F1.128C.Sample Abundance.F1.129N.Sample
## A0AVT1           5.444870           5.532511
## AOMZ66           4.735426           4.886217
##           Abundance.F1.129C.Sample Abundance.F1.130N.Sample
## A0AVT1           5.411776           5.353822
## AOMZ66           4.834322           4.645590
##           Abundance.F1.130C.Sample Abundance.F1.131.Sample
## A0AVT1           5.417480           5.350928
## AOMZ66           4.830086           4.634653
```

... and the associated feature data. Notice that there are many columns in the feature data. These are all the additional columns output from PD in addition to the quantification. They are all stored here in case they are required.

```
print(head(fData(total_protein_quant), 2))
```

```
##          Checked Confidence      Sequence
## A0AVT1   False           High    ACIGDTLCQK
## A0MZ66   False           High  ATQPETTEEVTDLK
##
##                                     Modifications Quality.PEP
## A0AVT1  2xTMT6plex [N-Term; K10]; 2xCarbamidomethyl [C2; C8] 0.000118596
## A0MZ66                                     2xTMT6plex [N-Term; K14] 0.000709357
##
##          Quality.q.value Number.of.Protein.Groups Number.of.Proteins
## A0AVT1   5.56241e-05                                1                1
## A0MZ66   0.000157133                                1                1
##
##          Number.of.PSMs Master.Protein.Accessions Number.of.Missed.Cleavages
## A0AVT1           3                                A0AVT1                0
## A0MZ66           3                                A0MZ66                0
##
##          Theo.MHplus.in.Da Quan.Info Amanda.Score.MS.Amanda
## A0AVT1   1623.85986925      Unique      330.291304522885
## A0MZ66   2020.08503667      Unique      342.226818892066
##
##          Confidence.MS.Amanda Search.Space.MS.Amanda
## A0AVT1           High                                1206.0
## A0MZ66           High                                1908.0
##
##          Percolator.q.Value.MS.Amanda Percolator.PEP.MS.Amanda
## A0AVT1           9.831e-05      0.00015900000000000002
## A0MZ66           0.0001625                                0.0004789
##
##          XCorr.Sequest.HT Confidence.Sequest.HT Search.Space.Sequest.HT
## A0AVT1   3.72960662841797                                High
## A0MZ66   3.85038113594055                                High
##
##          Percolator.q.Value.Sequest.HT Percolator.PEP.Sequest.HT
## A0AVT1           7.671e-05                                0.0001077
## A0MZ66           0.00014419999999999998                0.0003799
##
##          Ions.Score.Mascot Confidence.Mascot Search.Space.Mascot
## A0AVT1           71.39                                High
## A0MZ66           32.36                                High
##
##          Percolator.q.Value.Mascot Percolator.PEP.Mascot master_protein
## A0AVT1           8.725e-05 0.00010700000000000001      A0AVT1
## A0MZ66           0.0001791                                0.0007973      A0MZ66
##
##          protein_length protein_description peptide_start peptide_end
## A0AVT1           1052 sp|A0AVT1|UBA6_HUMAN                447        457
## A0MZ66           631 sp|A0MZ66|SHOT1_HUMAN                392        406
##
##          crap_protein associated_crap_protein unique
## A0AVT1           0                                0                1
## A0MZ66           0                                0                1
##
##                                     filename CV.Abundance.F1.126.Sample
## A0AVT1  Nocodazole_Total_PeptideGroups.txt                0.08912426
## A0MZ66  Nocodazole_Total_PeptideGroups.txt                0.14097226
##
##          CV.Abundance.F1.127N.Sample CV.Abundance.F1.127C.Sample
## A0AVT1           0.08717847                                0.1025999
## A0MZ66           0.10111405                                0.1889357
##
##          CV.Abundance.F1.128N.Sample CV.Abundance.F1.128C.Sample
## A0AVT1           0.12057913                                0.06179052
## A0MZ66           0.08956286                                0.13456950
##
##          CV.Abundance.F1.129N.Sample CV.Abundance.F1.129C.Sample
## A0AVT1           0.11457832                                0.04857286
## A0MZ66           0.04768461                                0.03199063
```

```
##          CV.Abandance.F1.130N.Sample CV.Abandance.F1.130C.Sample
## AOAVT1          0.05473802          0.1359917
## AOMZ66          0.13902945          0.2397172
##          CV.Abandance.F1.131.Sample
## AOAVT1          0.1632846
## AOMZ66          0.1071883
```

... and here is the phenotype data. As we can see, we have 3 replicates each of “M”, “G1” and “S” phase, plus an additional Control sample. For our purposes, we’re only going to be interested in the M and G1 phases so we can remove the other data. Both the total and RBP quantification objects have the exact same order

```
print(pData(total_protein_quant))
```

```
##          Sample_name Condition Replicate Type
## Abundance.F1.126.Sample Control_1 Control      1 Total
## Abundance.F1.127N.Sample      M_1      M      1 Total
## Abundance.F1.127C.Sample      M_2      M      2 Total
## Abundance.F1.128N.Sample      M_3      M      3 Total
## Abundance.F1.128C.Sample      G1_1      G1      1 Total
## Abundance.F1.129N.Sample      G1_2      G1      2 Total
## Abundance.F1.129C.Sample      G1_3      G1      3 Total
## Abundance.F1.130N.Sample      S_1      S      1 Total
## Abundance.F1.130C.Sample      S_2      S      2 Total
## Abundance.F1.131.Sample      S_3      S      3 Total
```

Below, we make a boolean from Conditions to subset the samples to those which are “M” or “G1”

```
# identify the samples we want to keep
samples_to_keep <- pData(total_protein_quant)$Condition %in% c("M", "G1")
print(samples_to_keep)
```

```
## [1] FALSE TRUE TRUE TRUE TRUE TRUE TRUE FALSE FALSE FALSE
```

```
total_protein_quant <- total_protein_quant[, samples_to_keep]
rbp_protein_quant <- rbp_protein_quant[, samples_to_keep]

print(pData(total_protein_quant))
```

```
##          Sample_name Condition Replicate Type
## Abundance.F1.127N.Sample      M_1      M      1 Total
## Abundance.F1.127C.Sample      M_2      M      2 Total
## Abundance.F1.128N.Sample      M_3      M      3 Total
## Abundance.F1.128C.Sample      G1_1      G1      1 Total
## Abundance.F1.129N.Sample      G1_2      G1      2 Total
## Abundance.F1.129C.Sample      G1_3      G1      3 Total
```

```
print(pData(rbp_protein_quant))
```

```
##          Sample_name Condition Replicate Type
## Abundance.F1.127N.Sample      M_1      M      1 OOPS
```

```
## Abundance.F1.127C.Sample      M_2      M      2 OOPS
## Abundance.F1.128N.Sample      M_3      M      3 OOPS
## Abundance.F1.128C.Sample      G1_1      G1      1 OOPS
## Abundance.F1.129N.Sample      G1_2      G1      2 OOPS
## Abundance.F1.129C.Sample      G1_3      G1      3 OOPS
```

To detect changes in RNA binding, we can only consider RBPs where we have also quantified the total protein. Below, we identify these cases by intersecting the rownames of each MSnSet (the protein names)

```
intersecting_proteins <- intersect(rownames(total_protein_quant), rownames(rbp_protein_quant))

cat(sprintf("Out of a total of %s RBPs quantified,\nwe have total protein quantification for %s proteins",
            length(rownames(rbp_protein_quant)),
            length(intersecting_proteins),
            round(100*length(intersecting_proteins)/length(rownames(rbp_protein_quant)), 2)))
```

```
## Out of a total of 2149 RBPs quantified,
## we have total protein quantification for 1916 proteins = 89.16 %
```

Below, we convert the MSnSet into a “tidy” format data.frame using `biobroom::tidy()`

```
total_exprs <- total_protein_quant[intersecting_proteins,] %>% # subset to intersecting proteins
  tidy(addPheno=TRUE) %>% # "tidy" the object, e.g make it into a tidy data format --> long
  mutate(intensity=value) %>% dplyr::select(-value) # rename the "value" column -> "intensity"
```

Top of the total protein expression data.frame. See how each intensity value now has it’s own row with the other columns describing the associated aspects of the intensity value, e.g the protein and experimental condition

```
print(head(total_exprs))
```

```
## # A tibble: 6 x 7
##   protein sample      Sample_name Condition Replicate Type  intensity
##   <fct>   <fct>      <chr>      <chr>      <chr>      <chr>      <dbl>
## 1 A0AVT1 Abundance.F1.127~ M_1      M          1      Total      5.38
## 2 A1L0T0 Abundance.F1.127~ M_1      M          1      Total      3.79
## 3 A1L390 Abundance.F1.127~ M_1      M          1      Total      4.90
## 4 A1X283 Abundance.F1.127~ M_1      M          1      Total      5.00
## 5 A5YKK6 Abundance.F1.127~ M_1      M          1      Total      5.22
## 6 A6NFI3 Abundance.F1.127~ M_1      M          1      Total      5.06
```

Question: Why doesn’t the MSnSet object store all the data in this long format?

Now we do the same for the RBP quantification and then concatenate the two data frames together.

```
oops_exprs <- rbp_protein_quant[intersecting_proteins,] %>%
  tidy(addPheno=TRUE) %>%
  mutate(intensity=value) %>% dplyr::select(-value)

combined_exprs <- rbind(total_exprs, oops_exprs)
```

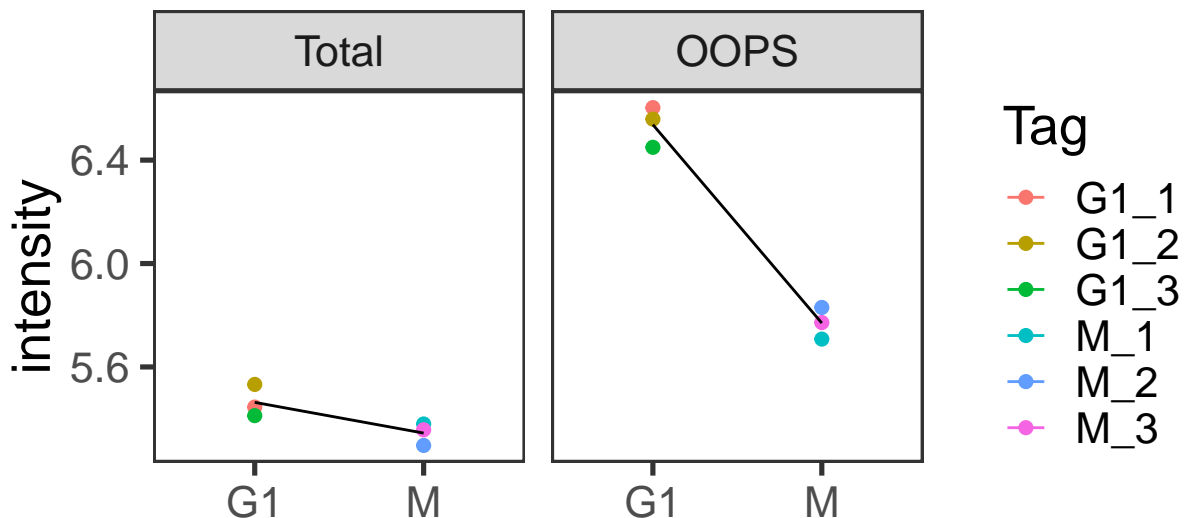
We want to tell R which is the order of the values in the condition and type columns so that the fold changes are in the expected direction, e.g positive = higher in G1 vs M.

```
combined_exprs$condition <- factor(combined_exprs$Condition, levels=c("M", "G1"))
combined_exprs$type <- factor(combined_exprs$Type, levels=c("Total", "OOPS"))
```

Now we model the protein intensity according to the models described in 1_simple_example_vd.Rmd. As an example, let's see the results from just applying the models to a single UniprotID.

```
combined_exprs %>% filter(protein == 'AOAVT1') %>%
  ggplot(aes(Condition, intensity, group=1, colour=factor(Sample_name))) +
  geom_point(size=2) +
  stat_summary(geom="line") +
  scale_colour_discrete(name="Tag") +
  xlab("") +
  facet_wrap(~type)
```

```
## No summary function supplied, defaulting to `mean_se()`
## No summary function supplied, defaulting to `mean_se()`
```



```
fit <- combined_exprs %>% filter(protein == 'AOAVT1') %>%
  lm(formula=intensity~Condition*Type)

print(summary(fit))
```

```
##
## Call:
## lm(formula = intensity ~ Condition * Type, data = .)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.087537 -0.048708  0.007263  0.041451  0.069459
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.53699    0.03620 180.590 9.89e-16 ***
## ConditionM       -0.76710    0.05119 -14.985 3.88e-07 ***
## TypeTotal        -1.07394    0.05119 -20.979 2.80e-08 ***
## ConditionM:TypeTotal  0.64836    0.07240   8.956 1.92e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0627 on 8 degrees of freedom
## Multiple R-squared:  0.988, Adjusted R-squared:  0.9835
## F-statistic: 219.6 on 3 and 8 DF, p-value: 5.07e-08
```

```
fit <- combined_exprs %>% filter(protein == 'AOAVT1') %>%
  lm(formula=intensity~Condition*Type+Sample_name)
print(summary(fit))
```

```
##
## Call:
## lm(formula = intensity ~ Condition * Type + Sample_name, data = .)
##
## Residuals:
##      1       2       3       4       5       6       7
## 0.048650 -0.054058  0.005408 -0.042081  0.023951  0.018130 -0.048650
##      8      9     10     11     12
## 0.054058 -0.005408  0.042081 -0.023951 -0.018130
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.560891    0.051612 127.119 2.30e-08 ***
## ConditionM       -0.783735    0.072991 -10.737 0.000426 ***
## TypeTotal        -1.073941    0.051612 -20.808 3.15e-05 ***
## Sample_nameG1_2     0.021609    0.063212   0.342 0.749662
## Sample_nameG1_3    -0.093305    0.063212  -1.476 0.213966
## Sample_nameM_1     -0.020733    0.063212  -0.328 0.759371
## Sample_nameM_2     -0.001057    0.063212  -0.017 0.987465
## Sample_nameM_3           NA         NA      NA      NA
## ConditionM:TypeTotal  0.648356    0.072991   8.883 0.000887 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06321 on 4 degrees of freedom
## Multiple R-squared:  0.9939, Adjusted R-squared:  0.9832
## F-statistic: 93.16 on 7 and 4 DF, p-value: 0.0002897
```

We can see that the model fits the data well (“Multiple R-squared: 0.9673, Adjusted R-squared: 0.955”). We can see that the interaction term that we’re interested in (for changes in RNA binding) significantly deviates from zero in both models.

Below, we make a function to run the linear models on a protein, select the best model and then return the required values from the model. When we run on the same protein as above, we can see that the best model is the one including the TMT tag as a co-variate.

```
testModels <- function(obj, coeff_of_interest="conditionG1:typeOOPS"){

  fit1 <- obj %>% lm(formula=intensity ~ condition + type + condition*type + sample)

  fit2 <- obj %>% lm(formula=intensity ~ condition + type + condition*type)

  if ( AIC(fit1) < AIC(fit2) ) {
    chosen_fit <- fit1
    fit_name <- "With_tag"
  } else {
    chosen_fit <- fit2
    fit_name <- "Without_tag"
  }

  fit_values <- c(round(coef(summary(chosen_fit))[coeff_of_interest,], 4),
                  round(summary(chosen_fit)$adj.r.squared, 4))

  names(fit_values) <- c("lm_fold_change", "lm_std_error", "lm_t_value", "lm_p_value", "lm_adj_R_squared")
  fit_values <- as.data.frame(t(fit_values), stringsAsFactors=FALSE)
  fit_values$fit <- fit_name

  return(fit_values)
}
```

```
combined_exprs %>% filter(protein == 'AOAVT1') %>% testModels()
```

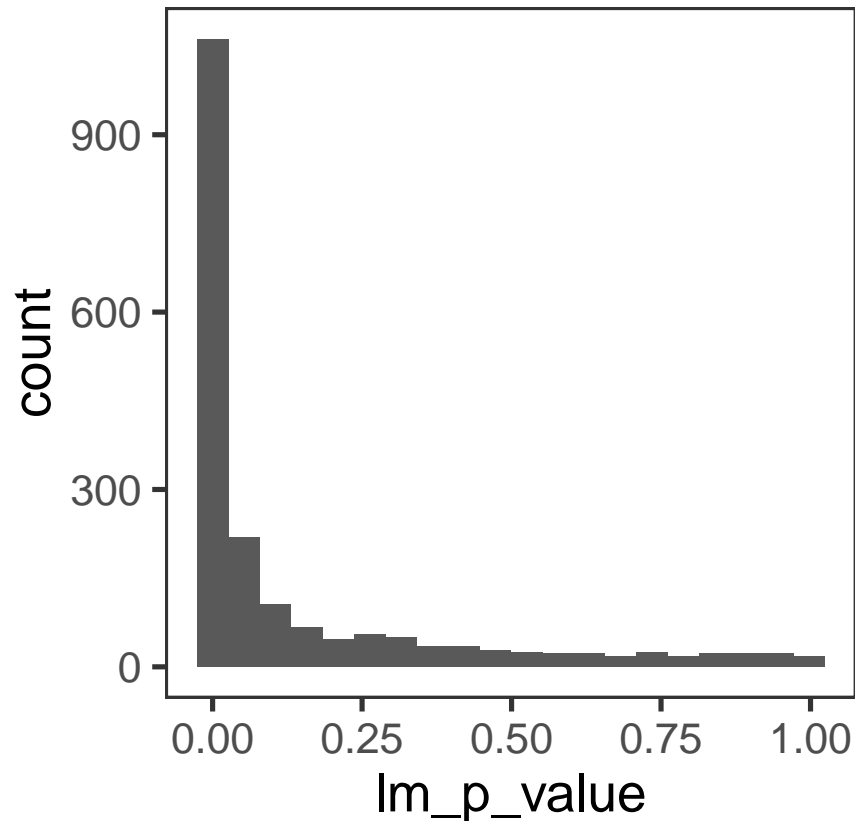
```
##   lm_fold_change lm_std_error lm_t_value lm_p_value lm_adj_R_squared
## 1          0.6484          0.073      8.8827      9e-04          0.9832
##      fit
## 1 With_tag
```

Below, we make a function to run the `testModels()` function on all proteins in turn using `dplyr`. We will use the standard Benjamini-Hochberg method to adjust p-values for the multiple tests we have conducted.

Below, we plot the p-values. Under the null hypothesis they should show an approximately uniform distribution. If there were a large number of proteins with a significant change in RNA binding, we would expect an additional “spike” with low p-values (<0.05). We see an approximately uniform distribution but with a slight skew towards low p-value. This may indicate the presence of changes in RNA binding but which we are insufficiently powered to detect, e.g low p-value but not significant low p-value.

```
plotP <- function(obj){
  p <- ggplot(obj, aes(lm_p_value)) + geom_histogram(bins=20)
  print(p)
}

plotP(M_G1)
```

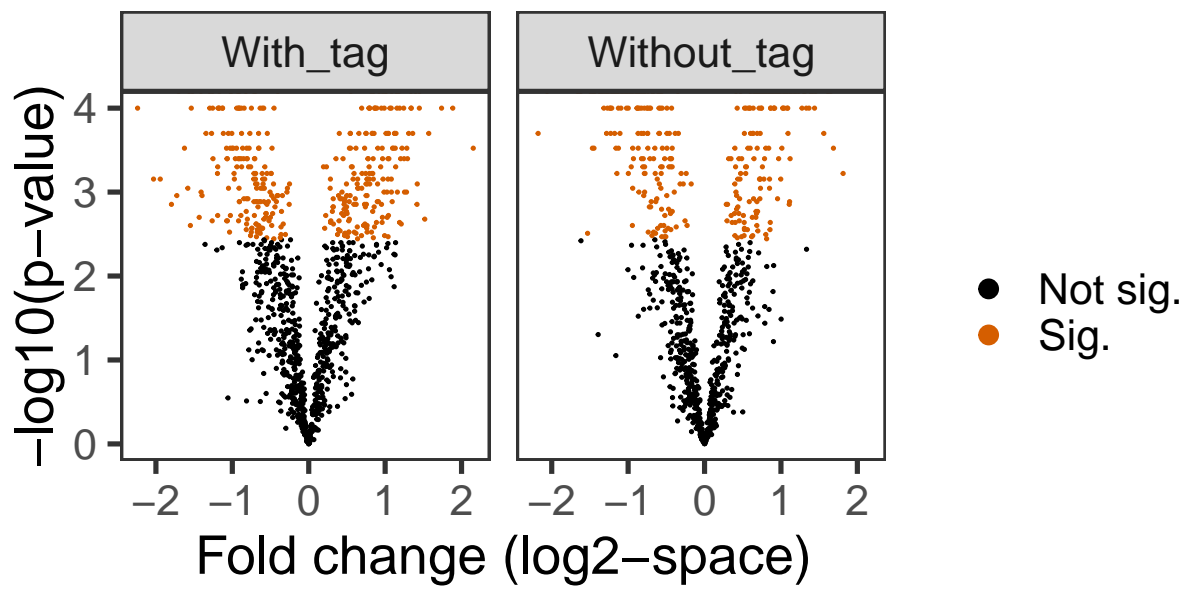
Questions:

- How many times was each of the two model formulas (with/without tags) used?
- How many significant changes in RNA binding were detected (You'll need to settle on a suitable FDR threshold) ?

So, we have detected a lot of proteins with a significant change in RNA binding!!

We can use volcano plots to take a look at the estimated fold changes and associated p-values

```
M_G1 %>%
  mutate(sig=ifelse(lm_BH<0.01, "Sig.", "Not sig.")) %>% # add "sig" column
  ggplot(aes(x=lm_fold_change, y=-log10(lm_p_value), colour=sig)) +
  geom_point(size=0.25) +
  scale_colour_manual(values=c("black", cbPalette[6]), name="") + # manually adjust colours
  facet_wrap(~fit) + # separate plots depending on model used
  xlab("Fold change (log2-space)") +
  ylab("-log10(p-value)") +
  guides(colour = guide_legend(override.aes = list(size = 3))) # manually set aesthetics for legend to
```



Finally, we save out the results for use in later notebooks

```
saveRDS(M_G1, "../results/M_G1_changes_in_RNA_binding_linear_model.rds")
```