

1 Reanalysis TMT Proteomics Data

At the center of the cogent critique of our manuscript was the questioned statistical validity of our previously described approach. Succinctly, the issue at question is whether or not the R package `edgeR` is an appropriate tool for analysis of protein mass spectrometry data.

Statistical inference in `edgeR` is built on a negative binomial (NB), generalized linear model (GLM) framework. Therefore, the data are assumed to be adequately described by a NB distribution parameterized by a dispersion parameter, ϕ .¹

Previously we used a customized workflow² to preprocess and normalize the data prior to performing statistical testing using `edgeR`'s flexible GLM framework. However, we failed to thoughtfully consider the overall adequacy of the NB framework for mass spectrometry data. Here we reconsider its appropriateness for our TMT proteomics dataset.

We evaluated the overall adequacy of the `edgeR` model by plotting the residual deviance of all proteins against their theoretical, normal quantiles in a quantile-quantile plot. **Figure 1** illustrates the overall lack of fit for the three dispersion models fit by `edgeR`. As an alternative to `edgeR` we considered `MSstatsTMT`, an extension of `MSstats` for analysis of TMT proteomics experiments.

`MSstatsTMT` utilizes a linear mixed-model framework. The strength of linear mixed models is their ability to account for complex sources of variation in an experimental design. In mixed models, one or more covariates are a categorical variable representing representing experimental or observational "units" in the data set (Bates 2010).

A TMT proteomics experiment consists of $m = 1 \dots M$ concatenations of isobaric TMT-labeled samples or **Mixtures**. Each TMT channel is dedicated to the analysis of $c = 1 \dots C$ individual biological or treatment **Conditions** prepared from one $b = 1 \dots B$ biological replicates or **Subjects**. A single mixture may be profiled in $t = 1 \dots T$ technical replicate mass spectrometry runs. Each protein is measured $M \times C \times B \times T$ times.

The full linear mixed model describing such an experiment is of the form:

$$Y_{mcbt} = \mu + Mixture_m + TechRepMixture_{t(m)} + Condition_c + Subject_{mcb} + \epsilon_{mcbt} \quad (1)$$

¹The dispersion parameter can take several forms. `edgeR` supports three dispersion models: 'common', 'trended', and 'tagwise'. However, when using `edgeR`'s robust quasi-likelihood test methods, only global (i.e. 'common' or 'trended') dispersion metrics are appropriate (see `?edgeR::glmQLFit`).

²The most important step in our normalization approach is IRS normalization. MS2 random sampling results in identification and quantification of proteins by different peptides in each MS experiment. To account for this source of variability, protein measurements are adjusted by a scaling factor such that the geometric mean of all internal reference standards are equal (Plubell et al., 2017). This is essential to account for the stochasticity of peptide quantification in MS experiments. Phillip Wilmarth's GitHub offers an excellent exploration of IRS normalization.

$$\begin{aligned}
Condition_c &= \sum_{C=1}^C \\
Mixture_m &\sim N(0, \sigma_M^2) \\
Subject_{mcb} &\sim N(0, \sigma_M^2) \\
TechRepMixture &\sim N(0, \sigma_T^2) \\
\epsilon_{mcbt} &\sim N(0, \sigma^2)
\end{aligned}$$

In an experiment with multiple mixtures and biological replicates, but no technical replication of mixture ($T = 1$) **MSstatsTMT** fits the model: Where **Mixture** is a mixed-effect and quantifies variation between TMT mixtures. **Condition** is a fixed effect (mean = 0) and in our experiment represents the interaction of terms **Genotype** and **BioFraction**. ϵ is a random effect representing both biological and technical variation, quantifying any remaining error.

If a term is not estimable, then it is removed from the model. In our experimental design, we made measurements from seven biological subcellular fractions (**BioFractions**) from each subject. Thus, we should include the term **Subject**, representing the 6 individual mice or subjects analyzed in our experiment. However, in our design **Mixture** is confounded with the term **Subject** – in each mixture we analyzed all **BioFractions** from a single Control and Mutant mouse. Thus we can choose to account for the effect of **Mixture** or **Subject**, but not both. Assuming **Mixture** contributes greater to the variance, we drop the term **Subject**, and the reduced model is equivalent to 2.

For experiments with multiple mixtures and biological replicates, the reduced model is then:

$$Y_{mcbt} = \mu + Mixture_m + Condition_c + \epsilon_{mcb} \quad (2)$$

Where **Condition** is the interaction of **Genotype** and **BioFraction**, the 14 combinations of 7 subcellular fractions measured in Control and Mutant mice.

Model based testing of differential abundance between pairs of conditions is assessed through contrast of conditioned means estimated by fitting the parameters of the model by REML to obtain $\hat{\beta}$, σ^2 and \hat{V} .

The degrees of freedom are determined by the Satterthwaite approximation[REF], and the T-statistic for the contrast is taken to be (`lmerTest ref`):

$$t = \frac{l^T * \hat{\beta}}{\sqrt{l * s^2 * \hat{V} * l^T}} \quad (3)$$

σ^2 is the error from **Equation 2**. $l^T = \sum_{C=1}^C = 0$ a vector specifying a contrast between positive and negative coefficients in the model.

Together the denominator $\sqrt{l * s^2 * \hat{V} * l^T}$ is the standard error of the contrast computed from the unscaled variance-covariance matrix, \hat{V} .

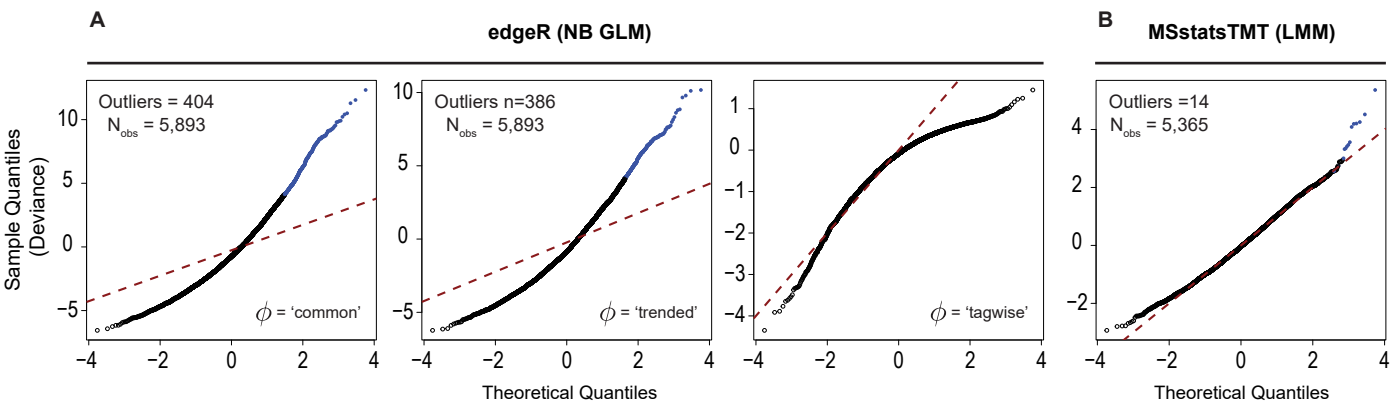


Figure 1: **Goodness-of-fit of edgeR (A), and MSstats (B) statistical approaches.** The overall adequacy of the linear models fit to the data were assessed by plotting the residual deviance for all proteins as a quantile-quantile plot (McCarthy *et al.*, (2012)). **(A)** The normalized protein data were fit with a NB GLM of the form: $\sim \text{Mixture} + \text{Condition}$. Where **Mixture** is a blocking factor that accounts for sources of variability between experiments. Protein-wise deviance statistics were transformed to normality and plotted against theoretical normal quantiles using `edgeR::gof`. **(B)** The normalized protein data were fit with a linear mixed-effects model (LMM) of the form: $\text{Abundance} \sim 0 + \text{Condition} + (1|\text{Mixture})$. Where **Mixture** indicates the random effect of **Mixture**. The residual deviance and degrees of freedom were extracted from the fitted models, z-score normalized, and plotted as in (A). Proteins with significantly poor fit are indicated as outliers in blue (Holm-adjusted P-value < 0.05).

	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13	C14	C15	C16
Mix1	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2
Mix2	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2
Mix3	WT-5K	WT-9K	WT-12K	WT-15K	WT-30K	WT-79K	WT-129K	QC1	MUT-5K	MUT-9K	MUT-12K	MUT-15K	MUT-30K	MUT-79K	MUT-129K	QC2

Figure 2: **Experimental Design**

```

> library(dplyr)
> library(lmerTest)
> library(SwipProteomics)

> data(msstats_prot)

> washc_prots = c('Q8C2E7', 'Q6PGL7', 'Q3UMB9', 'Q9CR27', 'Q8VDD8')

# fit the LMM

> fx <- Abundance ~ 0 + Genotype:BioFraction + (1|Mixture) + (1|Protein)

> fm <- lmer(fx, data=msstats_prot %>% filter(Protein %in% washc_prots))

> model_summary <- summary(fm, ddf='Satterthwaite')

```

Term	Estimate	SE	DF	Tvalue	Pvalue
Control:BioFractionF4	6.884	0.151	6.909	45.686	2.776e-09
Control:BioFractionF5	7.168	0.151	6.909	47.570	7.845e-10
Control:BioFractionF6	7.465	0.151	6.909	49.548	2.183e-09
Control:BioFractionF7	7.495	0.151	6.909	49.745	5.939e-10
Control:BioFractionF8	7.327	0.151	6.909	48.629	1.922e-09
Control:BioFractionF9	7.138	0.151	6.909	47.377	4.486e-10
Control:BioFractionF10	7.756	0.151	6.909	51.478	1.839e-09
Mutant:BioFractionF4	5.729	0.151	6.909	38.025	4.364e-10
Mutant:BioFractionF5	5.933	0.151	6.909	39.377	2.197e-09
Mutant:BioFractionF6	6.044	0.151	6.909	40.113	5.103e-10
Mutant:BioFractionF7	6.083	0.151	6.909	40.370	2.275e-09
Mutant:BioFractionF8	5.927	0.151	6.909	39.339	6.108e-10
Mutant:BioFractionF9	5.897	0.151	6.909	39.141	1.898e-09
Mutant:BioFractionF10	6.055	0.151	6.909	40.186	3.447e-10

Figure 3: Example: Fit lmer to Wash complex.

```

# Create a contrast to compare 'Control-Mutant'
> contrast <- lme4::fixef(fm)
> contrast[] <- 0
> contrast[grepl('Mutant', names(contrast))] <- +1/7 # Positive coeff
> contrast[grepl('Control', names(contrast))] <- -1/7 # Negative coeff

# Examine the results
> results <- lmerTestContrast(fm, contrast) %>%
>   mutate(Contrast = 'Mutate-Control') %>% unique() %>% knitr::kable()

```

Contrast	log2FC	percentControl	Pvalue	Tstatistic	SE	DF
Mutant-Control	-1.366434	0.3878488	0	-36.93673	0.0369939	190

Figure 4: Test for 'Mutant-Control' contrast for difference between means of WASH complex proteins.