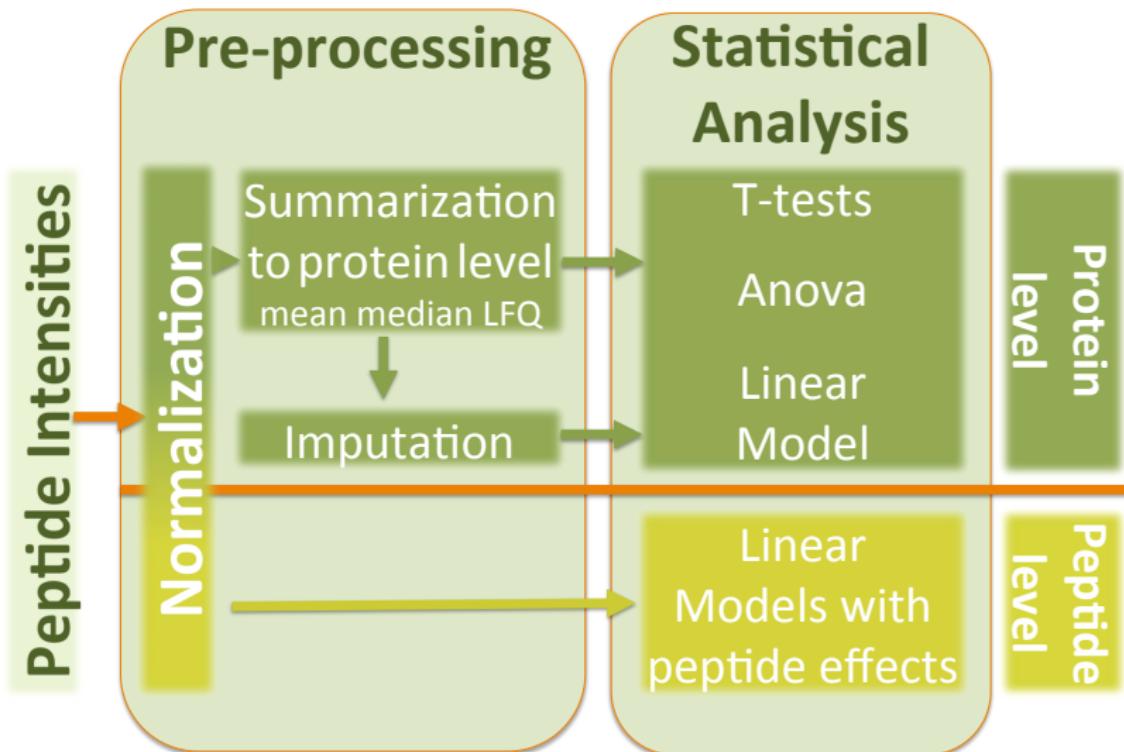


## Part II: Statistical Inference

Lieven Clement

Proteomics Data Analysis 2018, Gulbenkian Institute, May 28 -June 1  
2018.

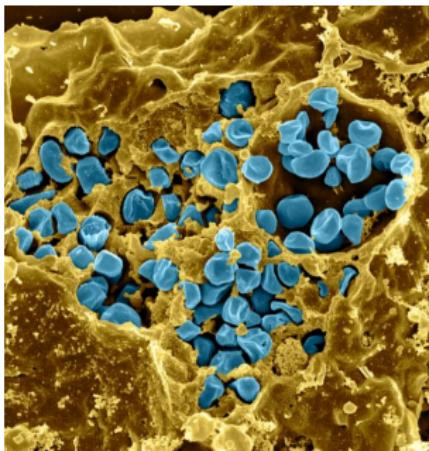
# Label-free Quantitative Proteomics Data Analysis Pipelines



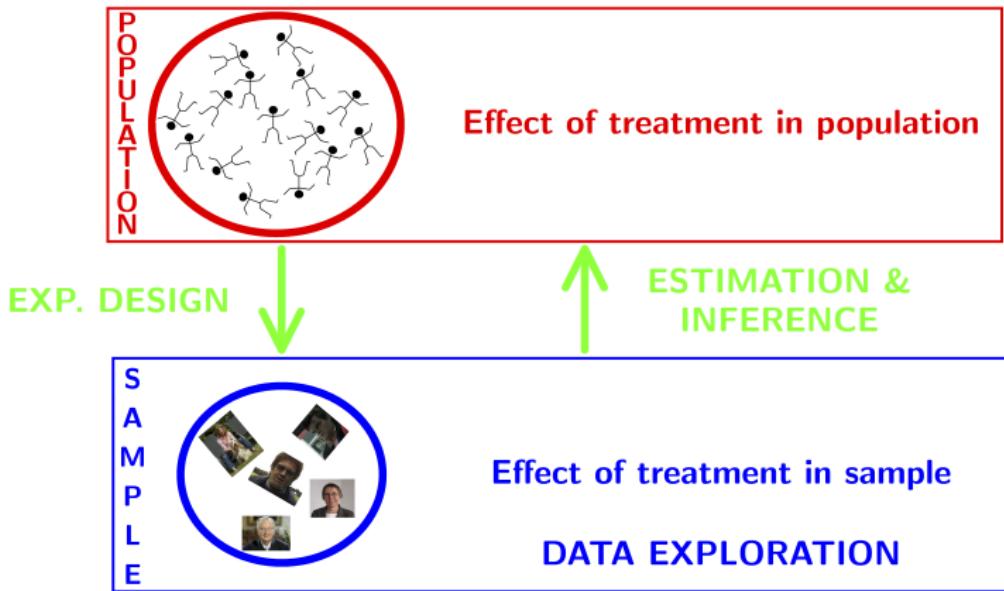
# Statistical Inference

- ① Francisella tularensis Example
- ② Hypothesis testing
- ③ Multiple testing
- ④ Moderated statistics
- ⑤ Experimental design
- ⑥ Peptide based models

# Francisella tularensis experiment



- Pathogen: causes tularemia
- Metabolic adaptation key for intracellular life cycle of pathogenic microorganisms.
- Upon entry into host cells quick phagosomal escape and active multiplication in cytosolic compartment.
- Francisella is auxotroph for several amino acids, including arginine.
- Inactivation of arginine transporter delayed bacterial phagosomal escape and intracellular multiplication.
- Experiment to assess difference in proteome using 3 WT vs 3 ArgP KO mutants

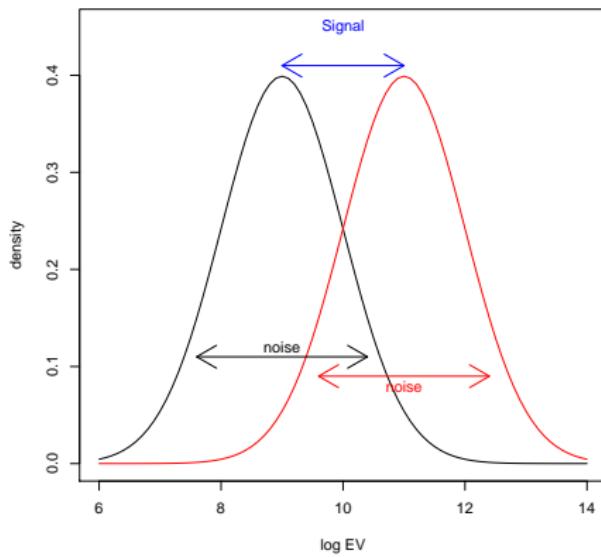


## Summarized data structure

- WT vs KO
- 3 vs 3 repeats
- 882 proteins

Protein	WT <sub>1</sub>	WT <sub>2</sub>	WT <sub>3</sub>	KO <sub>1</sub>	KO <sub>2</sub>	KO <sub>3</sub>
gi 118496616	29.83	29.77	29.91	29.70	29.86	29.80
gi 118496617	31.28	31.23	31.51	31.30	31.51	31.76
gi 118496635	32.39	32.27	32.24	32.25	32.14	32.22
gi 118496636	30.74	30.54	30.64	30.65	30.49	30.60
gi 118496637	29.56	29.35	29.56	29.30	29.24	29.14
gi 118498323	31.38	30.52	30.62	31.04	27.38	NA
:	:	:	:	:	:	:

# Hypothesis testing: a single protein



$$\Delta = \bar{z}_{p1} - \bar{z}_{p2}$$

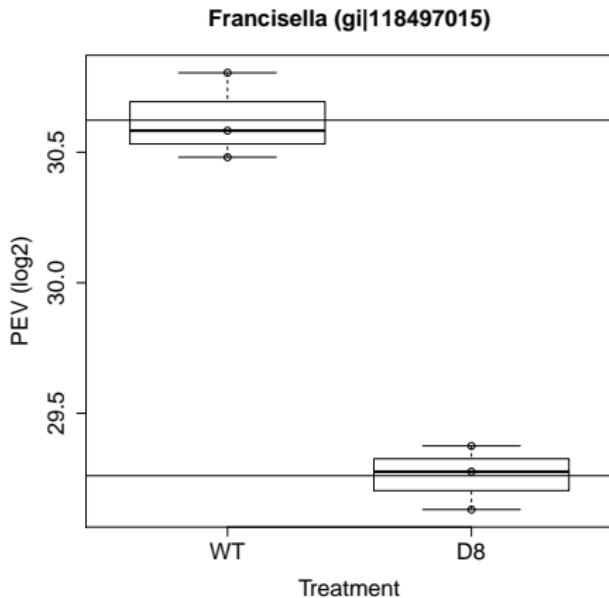
$$T_g = \frac{\Delta}{\text{se}_\Delta}$$

$$T_g = \frac{\overbrace{\text{signal}}^{\text{signal}}}{\overbrace{\text{Noise}}^{\text{Noise}}}$$

If we can assume equal variance in both treatment groups:

$$\text{se}_\Delta = \text{SD} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

# Hypothesis testing: a single protein



$$t = \frac{\log_2 \widehat{FC}}{\text{se}_{\log_2 \widehat{FC}}} = \frac{-1.4}{0.118} = -11.9$$

Is  $t = -11.9$  indicating that there is an effect?

How likely is it to observe  $t = -11.8$  when there is no effect of the argP KO on the protein expression?

# Null hypothesis and alternative hypothesis

- In general we start from **alternative hypothesis**  $H_A$ : we want to show an effect of the KO on a protein
  - On average the protein abundance in WT is different from that in KO

# Null hypothesis and alternative hypothesis

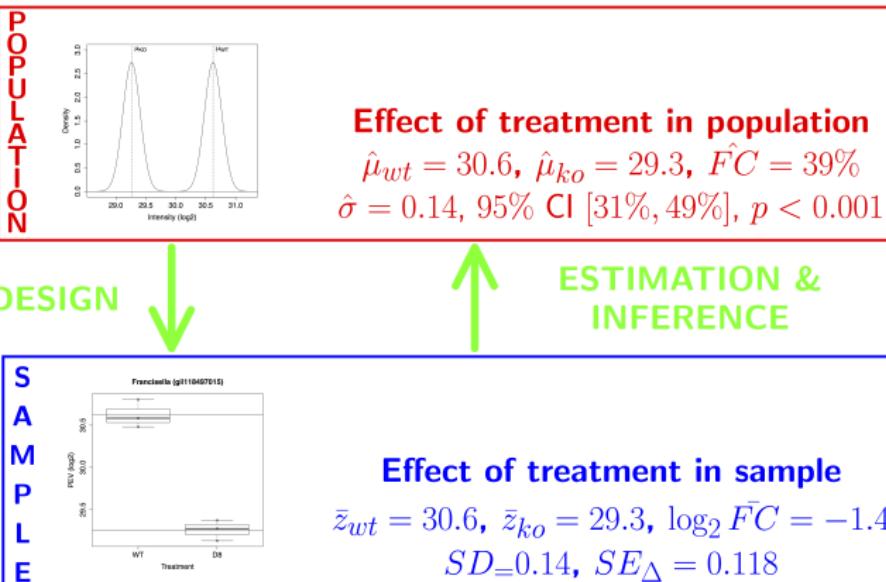
- In general we start from **alternative hypothesis**  $H_A$ : we want to show an effect of the KO on a protein
  - On average the protein abundance in WT is different from that in KO
- But, we will assess it by falsifying the opposite: **null hypothesis**  $H_0$ 
  - On average the protein abundance in WT is equal to that in KO

## Two Sample t-test

```
data: z by treat
t = -11.449, df = 4, p-value = 0.0003322
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-1.031371 -1.691774
sample estimates:
mean in group D8 mean in group WT
29.26094      30.62251
```

- How likely is it to observe an equal or more extreme effect than the one observed in the sample when the null hypothesis is true?
- When we make assumptions about the distribution of our test statistic we can quantify this probability: **p-value**. The p-value will only be calculated correctly if the underlying assumptions hold!
- When we repeat the experiment, the probability to observe a fold change more extreme than a 2.6 fold ( $\log_2 FC = -1.36$ ) down or up regulation by random chance (if  $H_0$  is true) is 3 out of 10.000.
- If the p-value is below a significance threshold  $\alpha$  we reject the null hypothesis. **We control the probability on a false positive result at the  $\alpha$ -level (type I error)**

# Hypothesis testing: a single protein



# Multiple hypothesis testing

# Problem of multiple hypothesis testing

- Consider testing DA for all  $m = 882$  proteins simultaneously
  - What if we assess each individual test at level  $\alpha$ ?
- Probability to have a false positive among all  $m$  simultaneous tests  $>>> \alpha = 0.05$

Suppose that 600 proteins are non-DA, then we could expect to discover on average  $600 \times 0.05 = 30$  false positive proteins. Hence, we are bound to call false positive proteins each time we run the experiment.

# FDR: False discovery rate

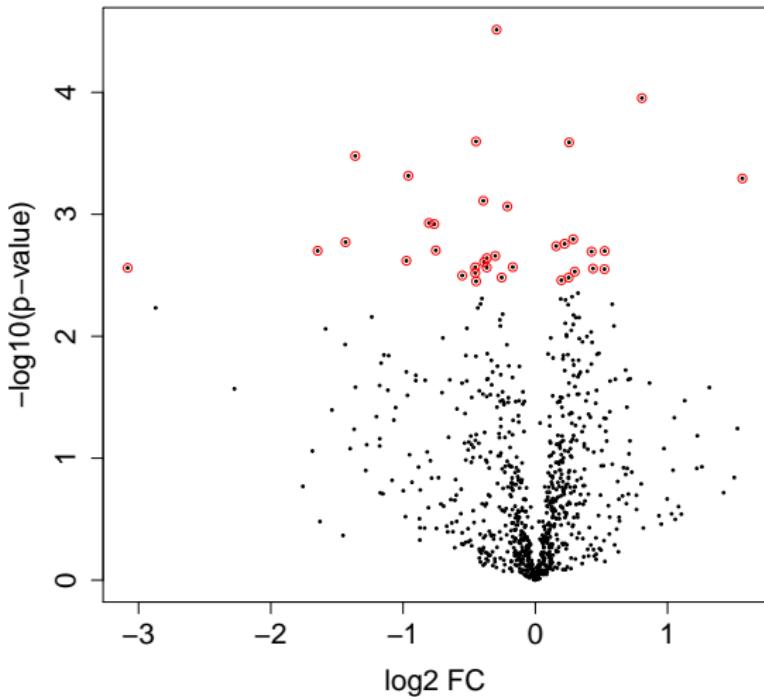
- FDR: Expected proportion of false positives on the total number of positives you return.
- An FDR of 1% means that on average we expect 1% false positive proteins in the list of proteins that are called significant.
- Defined by Benjamini and Hochberg in 1995

$$\text{FDR}(|t_{\text{thres}}|) = E \left[ \frac{FP}{FP + TP} \right] = \frac{\pi_0 Pr(|T| \geq t_{\text{thres}} | H_0)}{Pr(|T| \geq t_{\text{thres}})}$$

$$\text{FDR}_{\text{BH}}(|t_{\text{thres}}|) = \frac{1 \times p_{t_{\text{thres}}}}{\# |t_i| \leq t_{\text{thres}}} \frac{}{m}$$

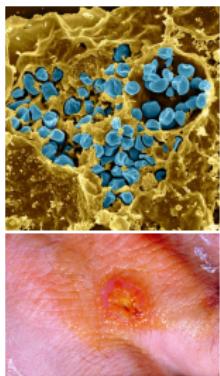
- FDR adjusted p-values can be calculated (e.g. Perseus, R, ...)

## Ordinary t-test

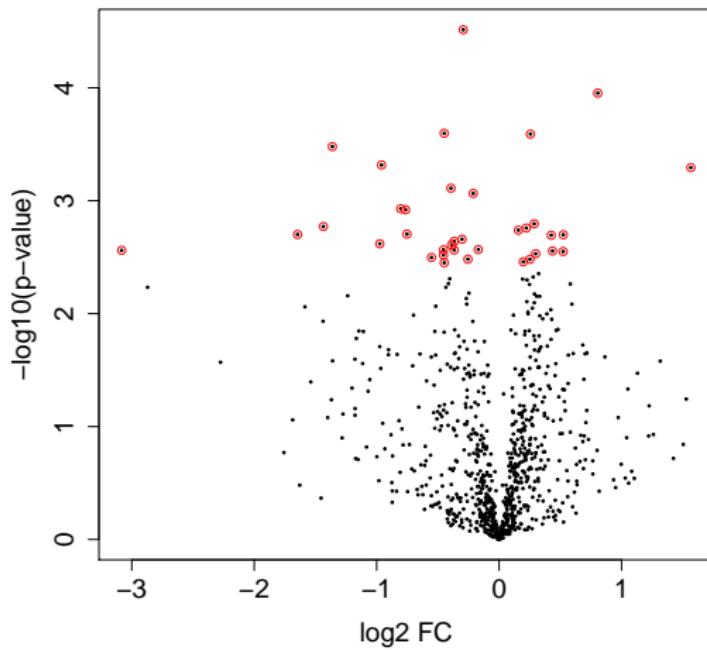


# Moderated Statistics

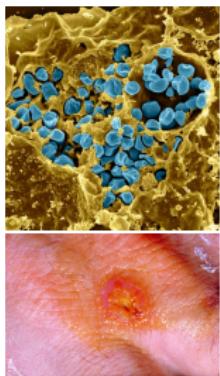
# Problems with ordinary t-test



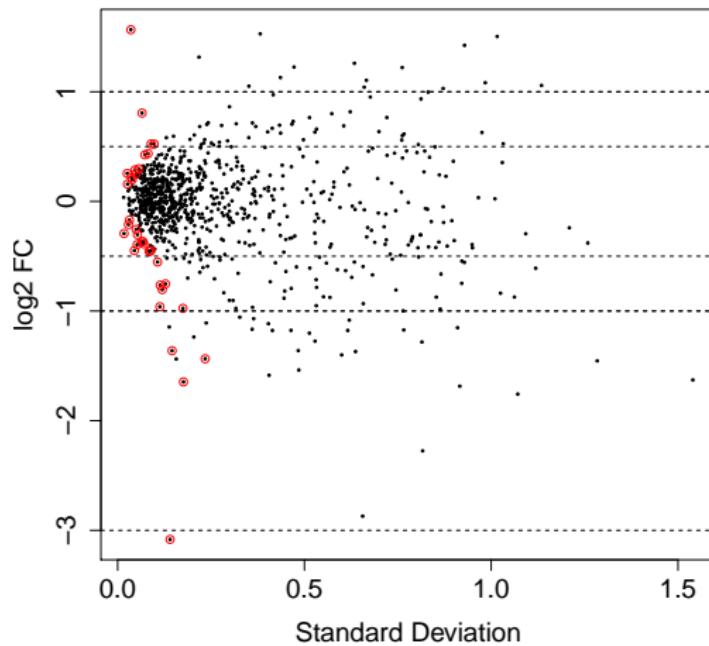
Ordinary t-test



# Problems with ordinary t-test

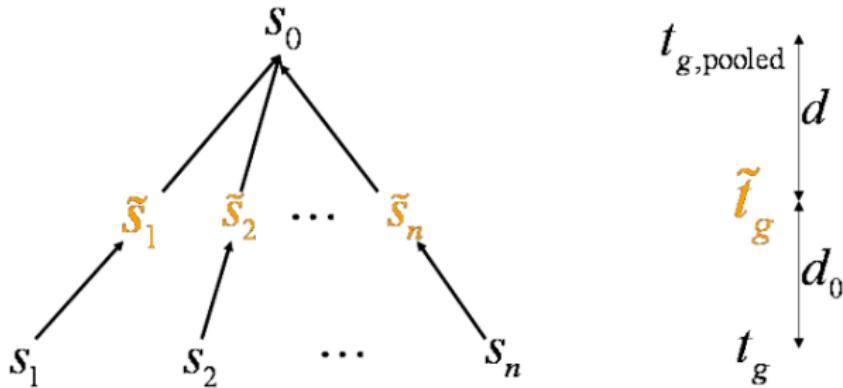


Original t-test



# Shrinkage of the variance and moderated t-statistics

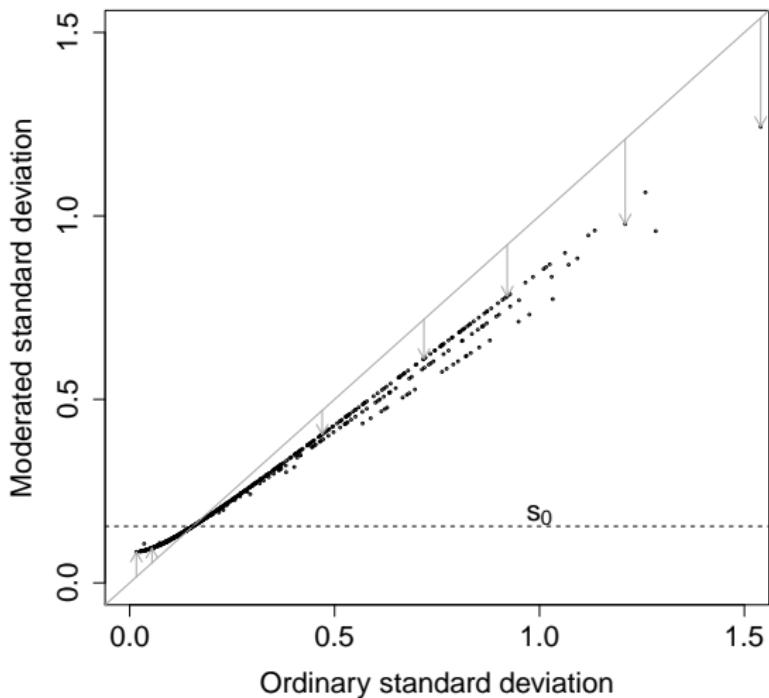
## Shrinkage of Standard Deviations



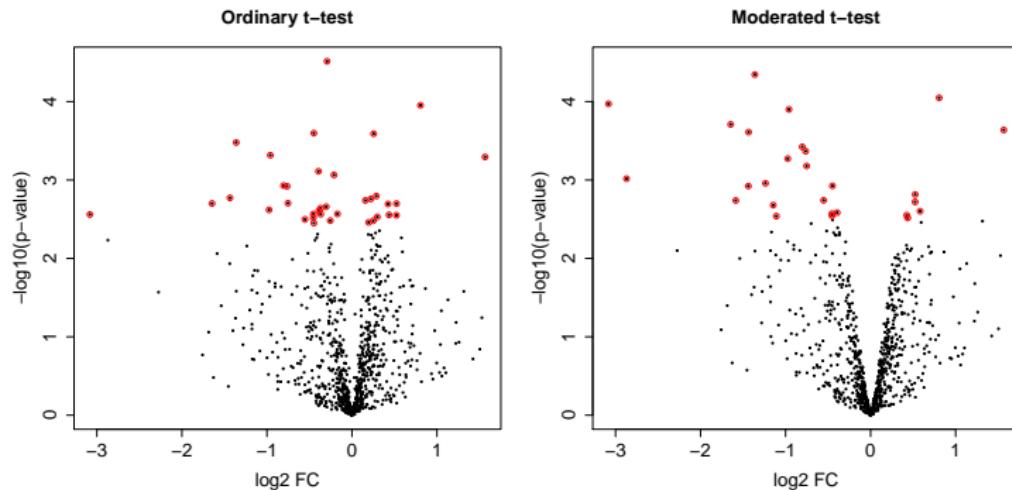
The data decides whether  $\tilde{t}_g$

should be closer to  $t_{g,\text{pooled}}$  or to  $t_g$

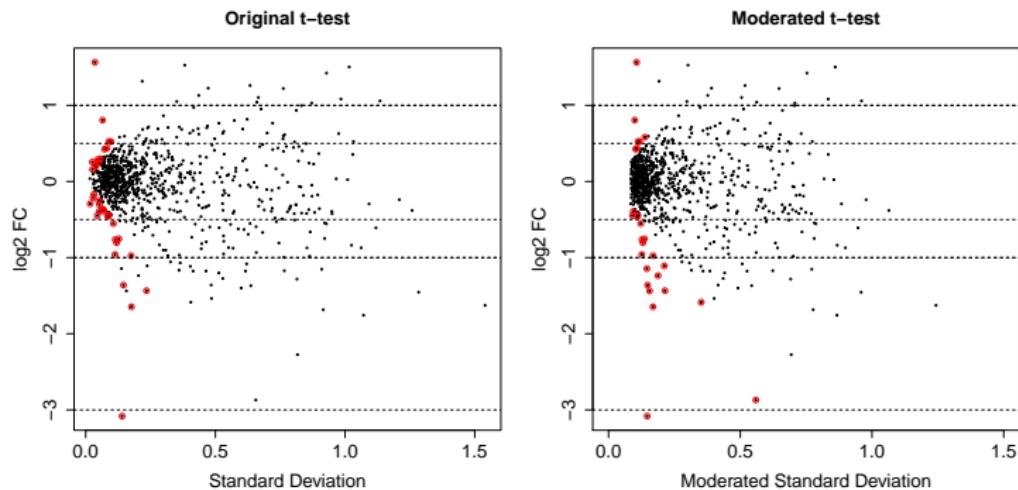
# Shrinkage of the variance with limma

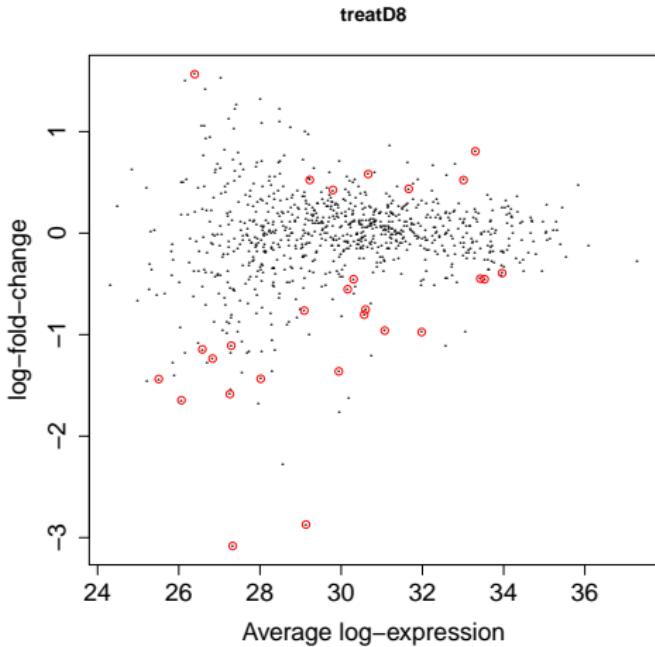


# Problems with ordinary t-test solved by moderated EB t-test



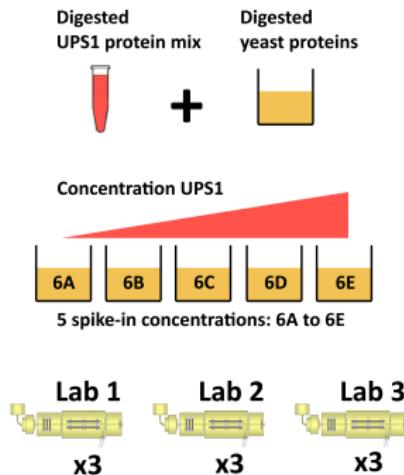
# Problems with ordinary t-test solved by moderated EB t-test





# Peptide-based models

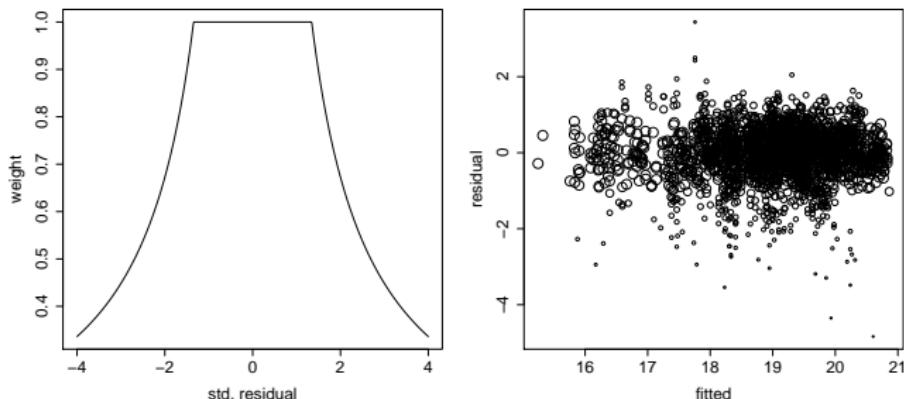
# Inference with Peptide Based Methods



- Protein by protein analysis of peptide level data with linear model  
 $y_{pept} \sim peptide + treatment + lab$
- Variance estimation in the literature:  
 protein-wise (LM) or via limma-style EB (LM-Sq).
- t-tests on model parameters

# Extension I: Robust estimation using observation weights (Ex I: LM-Sq-Rob)

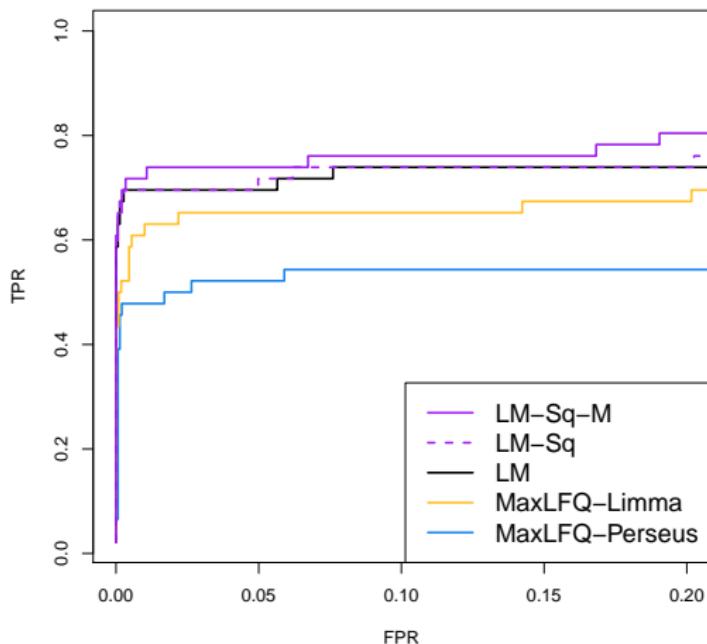
- Outlying peptide intensities: incorrect peptide identification, post-translational modifications, ...



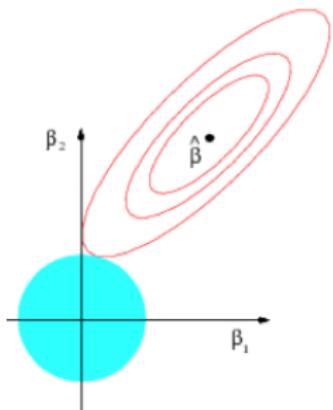
- Iteratively fit model with observation weights  $w(d_{ijp})$

$$\operatorname{argmin} \left[ \sum_{i=1}^n \sum_{j=p}^{P_j} w(d_{ijp}) \left( y_{ijp} - \mathbf{x}_i^T \boldsymbol{\beta}_j^{\text{treat}} - \beta_{jp}^{\text{pep}} \right)^2 \right]$$

# Method performance



## Extension II: Ridge regression

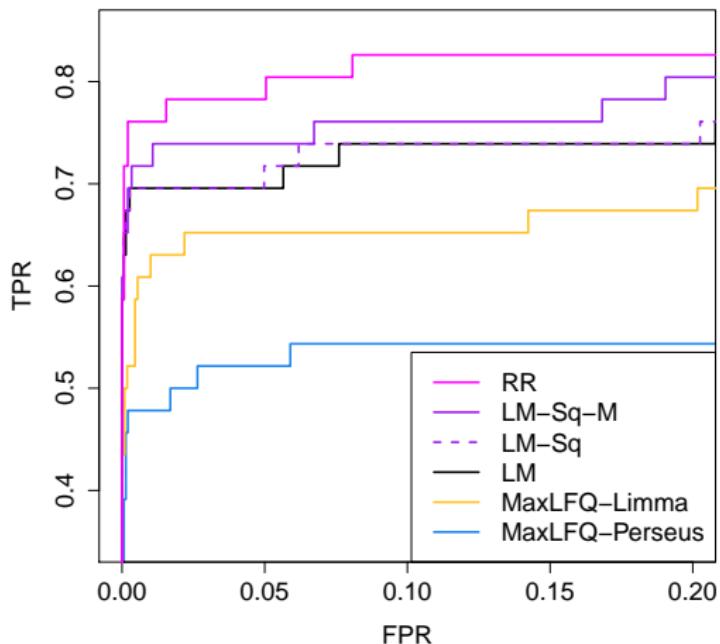


Parameters estimation via ridge regression,  
loss function:

$$\operatorname{argmin} \left[ \sum_{i=1}^n \sum_{j=1}^{P_j} w(d_{ijp}) \left( y_{ijp} - \mathbf{X}_i^T \boldsymbol{\beta}_j^{\text{treat}} - \beta_{jp}^{\text{pep}} \right)^2 + \lambda_j^{\text{treat}} \sum (\boldsymbol{\beta}_j^{\text{treat}})^2 + \lambda_j^{\text{pep}} \sum (\beta_{jp}^{\text{pep}})^2 \right]$$

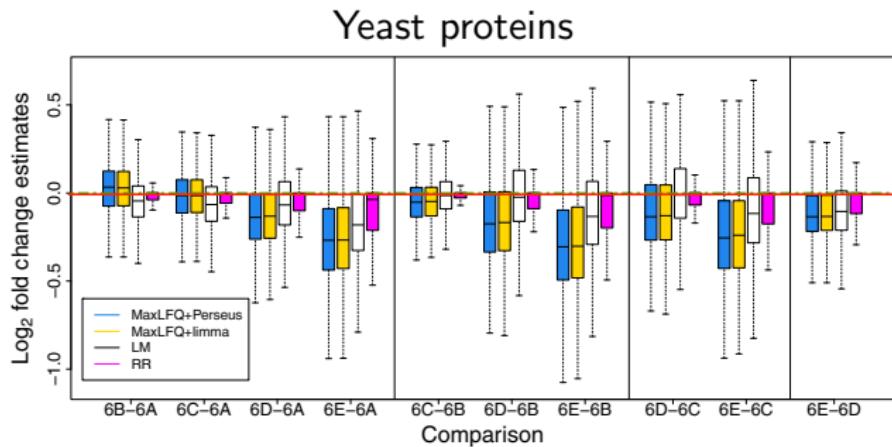
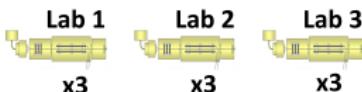
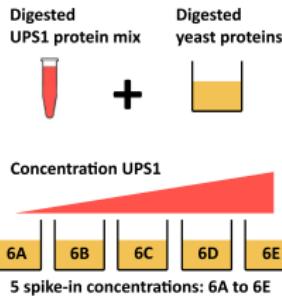
with

- $\lambda_{\text{treat}}$ : penalty term for regularization of parameters of interest
- $\lambda_{\text{pep}}$ : penalty term for regularization of peptide specific parameters



# Fold Change Estimates: Accuracy & Precision

## Study Design

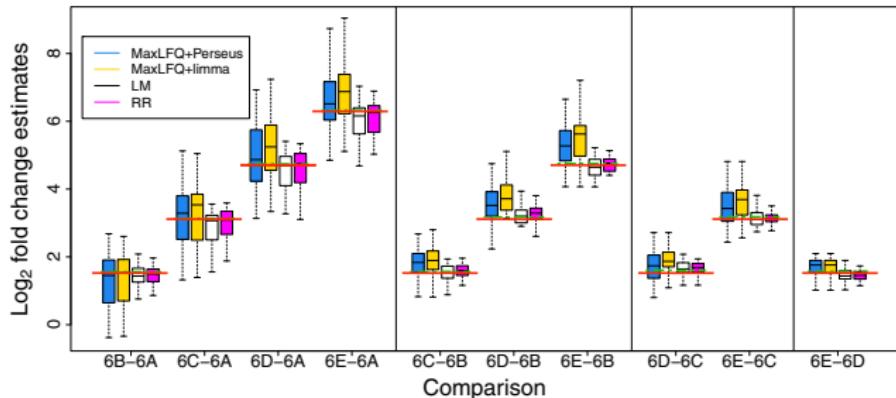
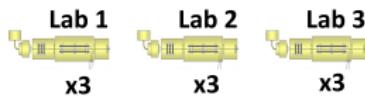
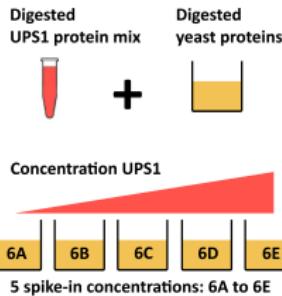


- Shrinkage: more precise and accurate FC estimates
  - Note, negative bias of the log<sub>2</sub> FC estimates as spike-in concentration increases
- Ionization suppression effects + Violation of normalization assumptions

# Fold Change Estimates: Accuracy & Precision

## Spiked UPS proteins

### Study Design



- MaxLFQ- Perseus and MaxLFQ-limma are always more biased and more variable
- Again MSqRob has a higher precision
- Shrinkage does not affect accuracy if there is evidence for DA!

# MSqRob

MSqRob for MaxQuant data v 0.8.0 Input Preprocessing Quantification

Select the grouping factor (mostly the "Proteins" column)

Proteins

Select additional annotation columns you want to keep

Protein names

Select fixed effects

genotype

Select random effects

Sequence run score

Save/Load options:

- Save the model
- Load existing models
- Don't save the model

Select the type of analysis

standard

Number of contrasts you want to test

1

Contrast 1

genotypeWT  
-1

genotypeKD  
1

Go



# GitHub

<https://github.com/statOmics/MSqRob>


- Goeminne, L., Gevaert, K. and Clement, L. (2016). Molecular and Cellular Proteomics, 15(2), 657-668
- Goeminne, L., Gevaert, K. and Clement, L. (2017). Journal of Proteomics, In Press.
- <http://dx.doi.org/10.1016/j.jprot.2017.04.004>