

## Model

## Abstract

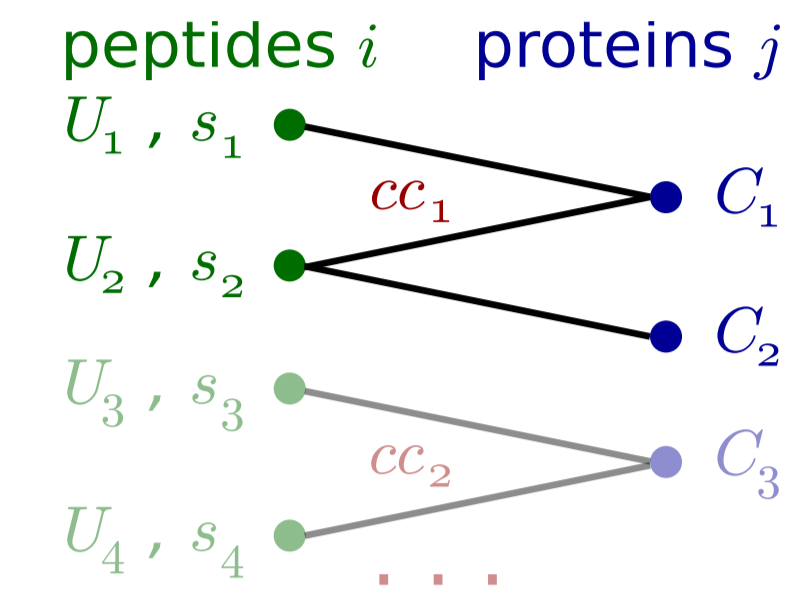
A major goal in proteomics is the comprehensive and accurate description of a proteome. Proteomics provides additional insights into biological systems that cannot be provided by genomic or transcriptomic approaches [1]. In particular, proteomics holds great promise for the identification and quantification of biomarkers capable of accurately predicting diseases already at a very early stage.

We propose a statistical approach to protein quantification based on shotgun experiments with four main advantages:

- Peptide intensities are modeled as random quantities, allowing to account for the uncertainty of these measurements.
- Our Markovian-type model for bipartite graphs ensures transparent propagation of the uncertainties and reproducible results.
- The problem of peptides mapping to several protein sequences (often neglected in other models) is addressed automatically according to our statistical model.
- Peptides with outlying intensity values can be assessed and classified as either regular data points or “true” outliers.

The application of our model is illustrated on two data sets and the performance is compared to another common approach for protein quantification [2].

## Notation



$U_i$  : intensity score (given)  
 $s_i$  : identification score (given)  
 $C_j$  : concentration (unknown)  
 → latent variable  
 $cc_r$  : connected component with  
 •  $n_r$  peptides  
 •  $m_r$  proteins

We define  $\underline{U}^{(r)}$  as the vector of intensities of all peptides in the connected component  $r$ . By writing  $\underline{\alpha}$  we mean  $\alpha(1, \dots, 1)^T$ .

Furthermore, we use a “distance” matrix  $D$  with

- $D_{ii}$  = number of proteins having a common edge with peptide  $i$
- $D_{ik}$  = number of proteins having a common edge with peptides  $i$  and  $k$

## Markovian-type assumptions

- peptides belonging to the same connected component are independent given their matching proteins  
 → dependencies among peptides are exclusively due to their common proteins
- only neighboring proteins matter in the conditional distribution for the peptides (see also [3])

## Model

We propose the following model for the peptide intensities:

$$U_i = \alpha + s_i \beta \sum_{j \in Ne(i)} C_j + \epsilon_i, \text{ with}$$

$$C_1, C_2, \dots, C_m \stackrel{i.i.d.}{\sim} \mathcal{N}(\mu, 1)$$

$$\epsilon_1, \epsilon_2, \dots, \epsilon_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \tau^2)$$

The elements of the covariance matrix of  $U$  are then given by

$$\text{Cov}(U_i, U_k) = \Sigma_{ik} = \begin{cases} s_i s_k \beta^2 D_{ik} & \text{for } i \neq k \\ s_i^2 \beta^2 D_{ii} + \tau^2 & \text{for } i = k \end{cases}$$

and the covariance between  $C_j$  and  $U_i$  is

$$\text{Cov}(C_j, U_i) = \Gamma_{C_j U_i} = \begin{cases} 0 & \text{if there is no edge between } i \text{ and } j \\ s_i \beta & \text{if there is an edge between } i \text{ and } j \end{cases}$$

## Predicting protein concentrations

For a connected component  $r$ :

$$\mathbb{E}[C_j | \underline{U}^{(r)}] = \mu + (\underline{U}^{(r)} - \underline{\alpha} - \underline{s}^{(r)} \beta \mu \text{diag}(D^{(r)}))^\top \Sigma_{\underline{U}^{(r)}}^{-1} \Gamma_{C_j \underline{U}^{(r)}}$$

## Parameter estimation

## Maximum likelihood estimation (MLE)

$$\underline{U}^{(r)} \sim \mathcal{N}_n(\underline{m}, \Sigma_{\underline{U}^{(r)}}) \text{ with } \underline{m} = \underline{\alpha} + \underline{s}^{(r)} \beta \mu \text{diag}(D^{(r)})$$

$$f(\underline{U}^{(r)}; \alpha, \beta, \mu, \tau^2) = |2\pi \Sigma_{\underline{U}^{(r)}}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\underline{U}^{(r)} - \underline{m})^\top \Sigma_{\underline{U}^{(r)}}^{-1} (\underline{U}^{(r)} - \underline{m})\right)$$

Minimize negative log-likelihood with respect to  $\alpha, \beta, \mu$  and  $\tau^2$

$$-\sum_{r=1}^R \log(f(\underline{U}^{(r)}; \alpha, \beta, \mu, \tau^2))$$

## Least squares approach (LSA)

Estimate  $\alpha$  and  $\beta \mu$  by fitting  $U \sim s \cdot \text{diag}(D)$ . Use sample covariance matrix of  $U$  to estimate  $\beta$  and  $\tau$ .

Off-diagonal elements of  $\hat{\Sigma}_{\underline{U}^{(r)}}$  allow to estimate  $\beta$ :

$$\sum_{r=1}^R \sum_{\substack{i \neq k \\ i, k \in cc_r}} \left( (\hat{\Sigma}_{\underline{U}^{(r)}})_{ik} - s_i s_k D_{ik} \beta^2 \right)^2 \stackrel{!}{=} \text{minimize w.r.t. } \beta^2$$

Diagonal elements of  $\hat{\Sigma}_{\underline{U}^{(r)}}$  and  $\hat{\beta}$  yield an estimate for  $\tau^2$ :

$$\sum_{r=1}^R \sum_{i=1}^{n_r} \left( (\hat{\Sigma}_{\underline{U}^{(r)}})_{ii} - s_i^2 \hat{\beta}^2 D_{ii} - \tau^2 \right)^2 \stackrel{!}{=} \text{minimize w.r.t. } \tau^2$$

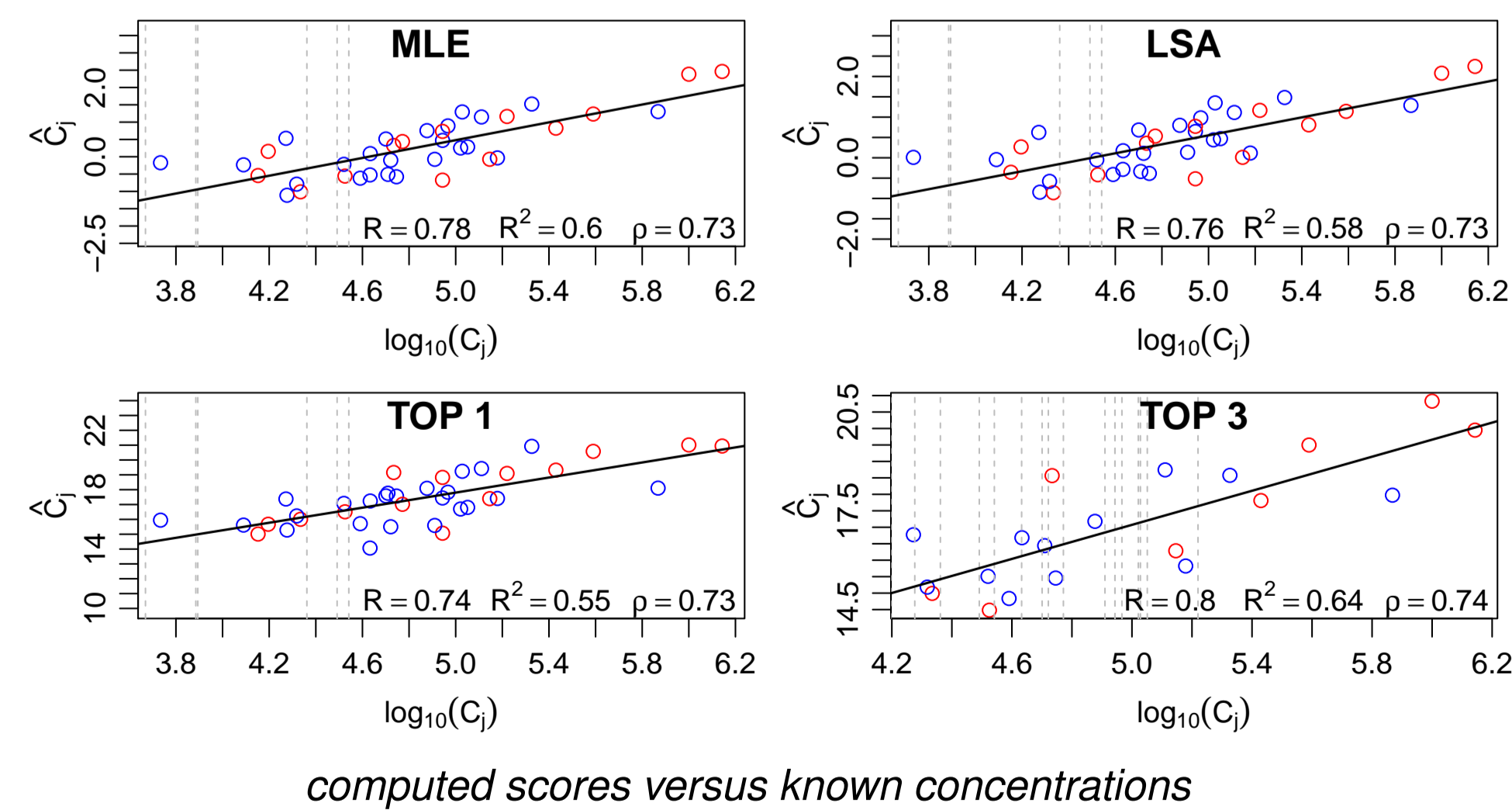
## Results

## Human shotgun data [4]

## “Absolute” protein quantification

Shotgun experiment (18 OGE fractions) on Human samples under various conditions. Performances of different quantification methods are compared on 44 experimentally quantified proteins (using AQUA peptides [5]) in the control sample.

Red symbols correspond to proteins identified (partly) with shared peptides. Grey lines indicate ground truth proteins that could not be quantified.



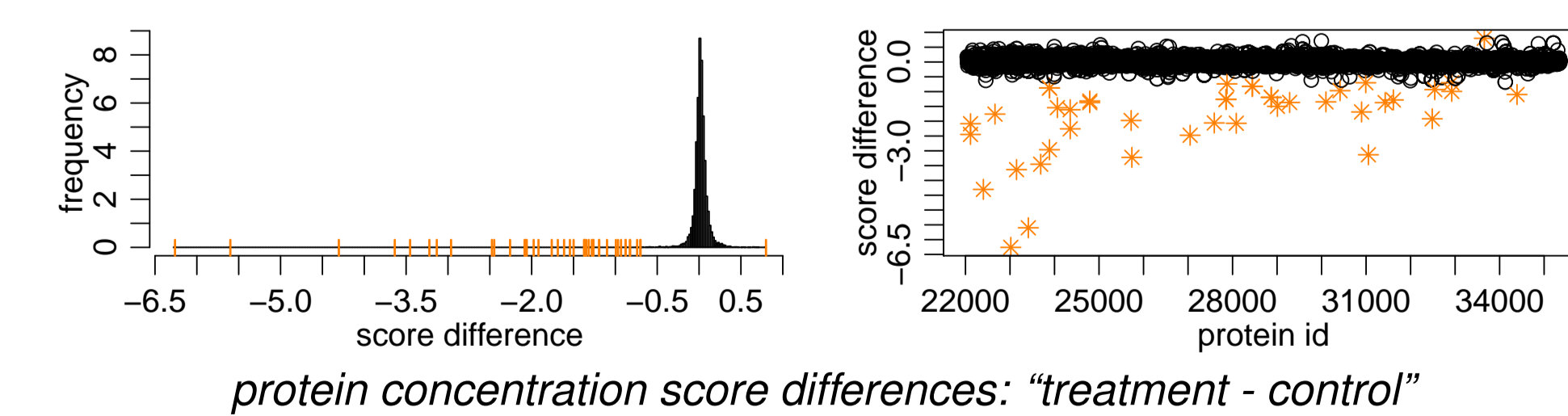
## SILAC-labeled human shotgun data [6]

Comparison between a control sample and irradiated (SILAC-labeled) human cells. The data sets contain much more shared peptides than unique ones. The bipartite graph holds:

- 22040 peptide and 13315 protein sequences
- 16728 of the peptides match to several proteins (shared)
- 12743 proteins match to at least one shared peptide
- 12018 proteins have no unique peptide evidence
- only 1297 proteins can be quantified with TOP1 (560 with TOP3)

## Relative protein quantification

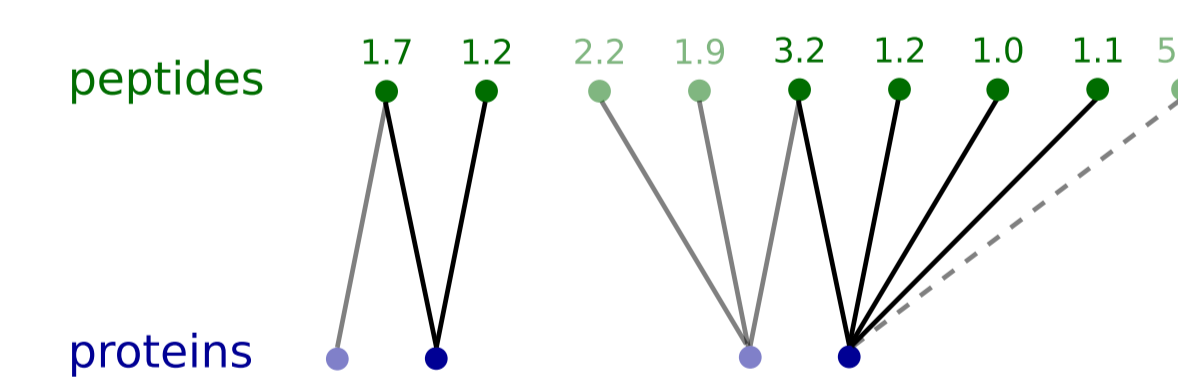
Proteins are quantified separately for the control and treatment sample, respectively. Score differences are used to find the proteins with the most important concentration changes between the two samples. Proteins with particularly high score differences are shown in orange.



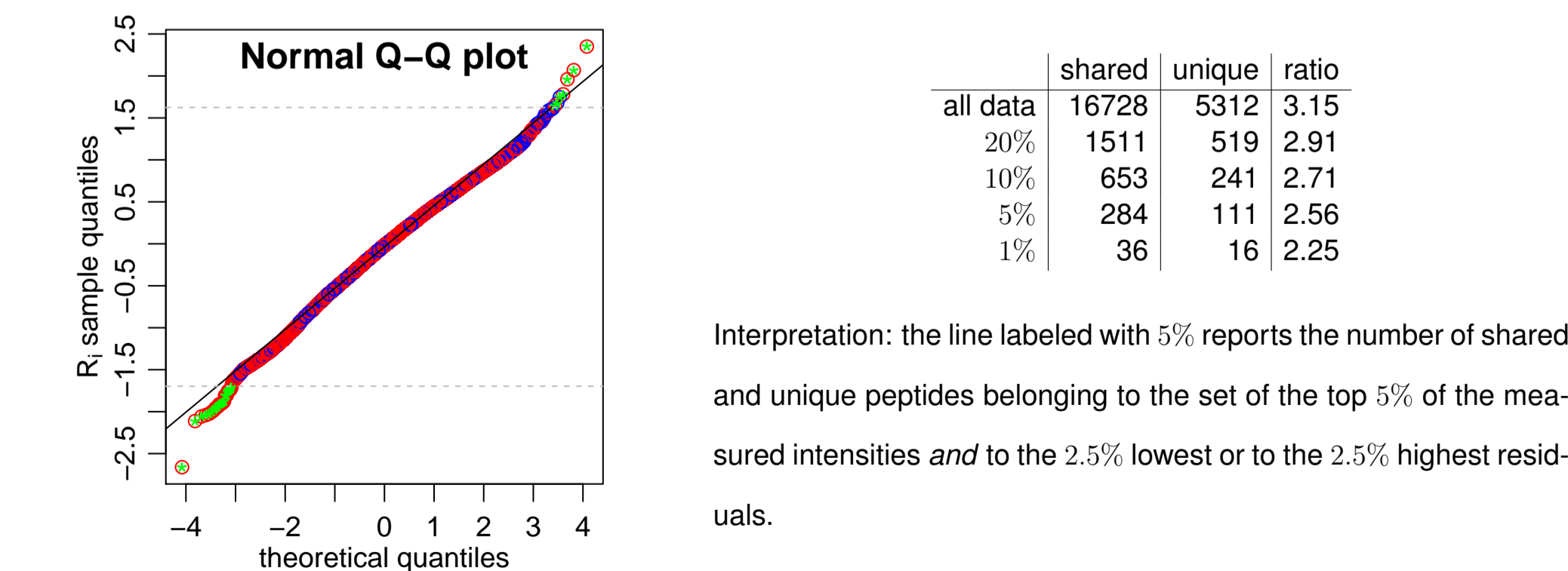
## Peptide intensity outlier detection

There are various reasons to encounter outliers in the measured peptide intensities:

- measurement errors
- shared peptides
- modifications
- missed cleavages
- incomplete database



Predicting peptide intensities and analyzing the residuals can potentially lead to further insight about the data. Monitoring peptides with large intensity measurements which could not be explained by the model (high residuals) will for example help to assess how useful the new approach is when dealing with shared peptides.



## Conclusion &amp; outlook

We provide a simple model with clearly stated assumptions for protein quantification.

- performance similar to other approaches
- when working on data sets with shared peptides our model allows to quantify all proteins
- potential to get further insight on the peptide level (assessing intensity measurements)
- flexible (peptide scores, input type, experiment type,...)

## Implementation

Our model is implemented in R [7]. The presented results are obtained with the following version:

- R version 2.15.0 RC (2012-03-25 r58832), x86\_64-unknown-linux-gnu
- Base packages: base, data sets, graphics, grDevices, methods, splines, stats, utils
- Other packages: MASS 7.3-17, regr0 1.0-2, sfsmisc 1.0-20, xtable 1.7-0
- Loaded via a namespace (and not attached): tools 2.15.0

## References

[1] Ruedi Aebersold and Matthias Mann. Mass spectrometry-based proteomics. *Nature*, 422(6928):198–207, 2003.  
 [2] Jeffrey C. Silva, Marc V. Gorenstein, Guo-Zhong Li, Johannes P. C. Vissers, and Scott J. Geromanos. Absolute quantification of proteins by LCMSE: A virtue of parallel ms acquisition. *Molecular & Cellular Proteomics*, 5:144–156, 2006.

[3] Sarah Gerster, Ermir Qeli, Christian H. Ahrens, and Peter Bühlmann. Protein and gene model inference based on statistical modeling in k-par tite graphs. *Proceedings of the National Academy of Sciences*, 107(27):12101–12106, 2010.

[4] Martin Beck, Alexander Schmidt, Johan Malmstroem, Manfred Claassen, Alessandro Ori, Anna Szymborska,

Franz Herzog, Oliver Rinner, Jan Ellenberg, and Ruedi Aebersold. The quantitative proteome of a human cell line. *Molecular Systems Biology*, 7(549), 2011.

[5] Scott A. Gerber, John Rush, Olaf Stemman, Marc W. Kirschner, and Steven P. Gygi. Absolute quantification of proteins and phosphoproteins from cell lysates by tandem MS. *Proceedings of the National Academy of Sciences*,

100(12):6940–6945, 2003.

[6] Mariette Matondo. Silac-labeled human shotgun data. personal communication.

[7] R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2011. ISBN 3-900051-07-0.