

Part I: Normalization & Summarization

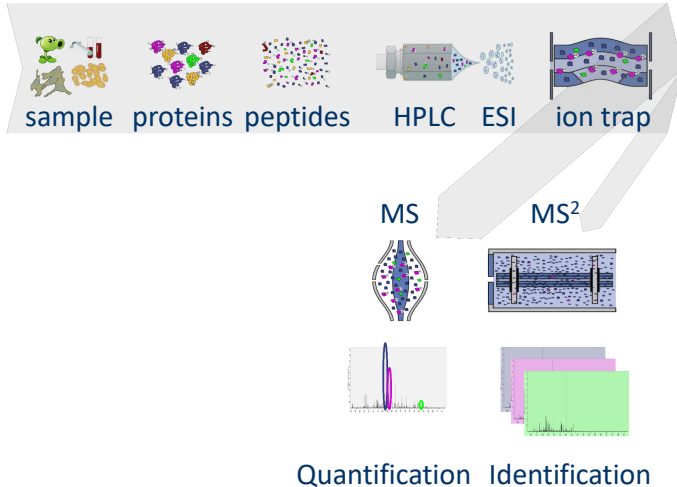
Lieven Clement

Proteomics Data Analysis Shortcourse

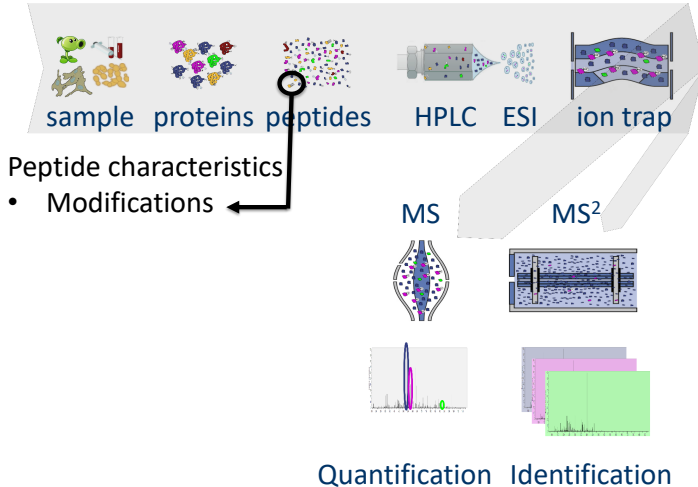
Outline

- 1 Introduction
 - 1 Label free MS based Quantitative Proteomics Workflow and Challenges
- 2 Preprocessing
 - 1 Filtering
 - 2 Log transformation
 - 3 Normalization
 - 4 Summarization

Challenges in Label Free Quantitative Proteomics



Challenges in Label Free Quantitative Proteomics

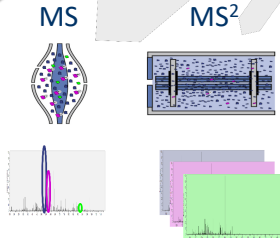


Challenges in Label Free Quantitative Proteomics



Peptide characteristics

- Modifications
- Ionisation efficiency
 - Outliers
 - Huge variability



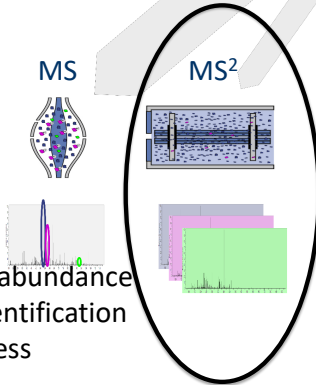
Quantification Identification

Challenges in Label Free Quantitative Proteomics



Peptide characteristics

- Modifications
- Ionisation efficiency
 - Outliers
 - Huge variability
- MS² selection on peptide abundance
 - Context dependent Identification
 - Non-random missingness



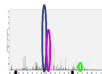
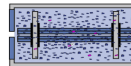
Challenges in Label Free Quantitative Proteomics



Peptide characteristics

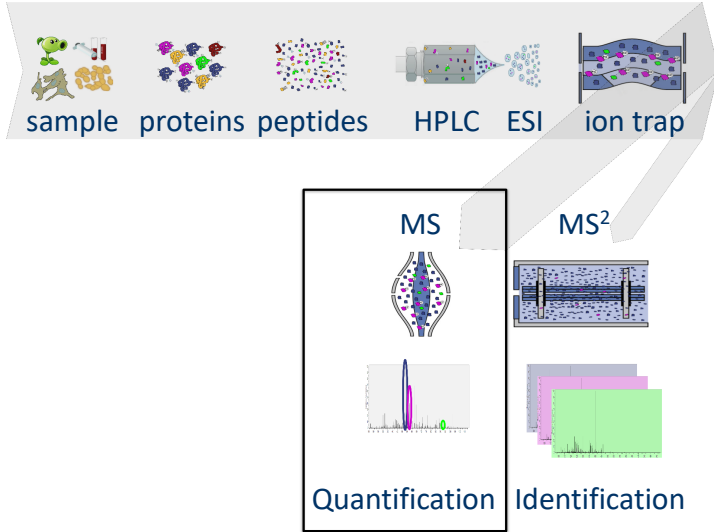
- Modifications
- Ionisation efficiency
 - Outliers
 - Huge variability
- MS² selection on peptide abundance
 - Context dependent Identification
 - Non-random missingness

MS

MS²

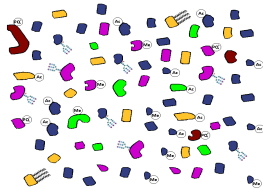
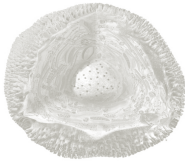
Unbalanced peptides identifications across samples and messy data

Challenges in Label Free MS-based Quantitative proteomics



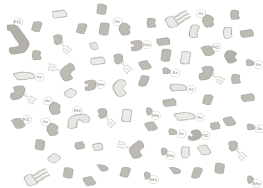
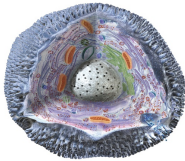
Challenges in Label Free MS-based Quantitative proteomics

MS-based proteomics returns **peptides**:
pieces of proteins

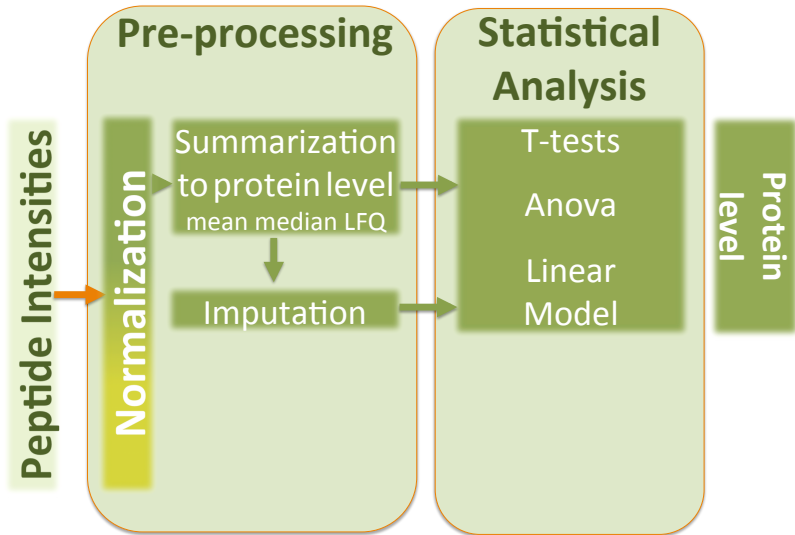


Challenges in Label Free MS-based Quantitative proteomics

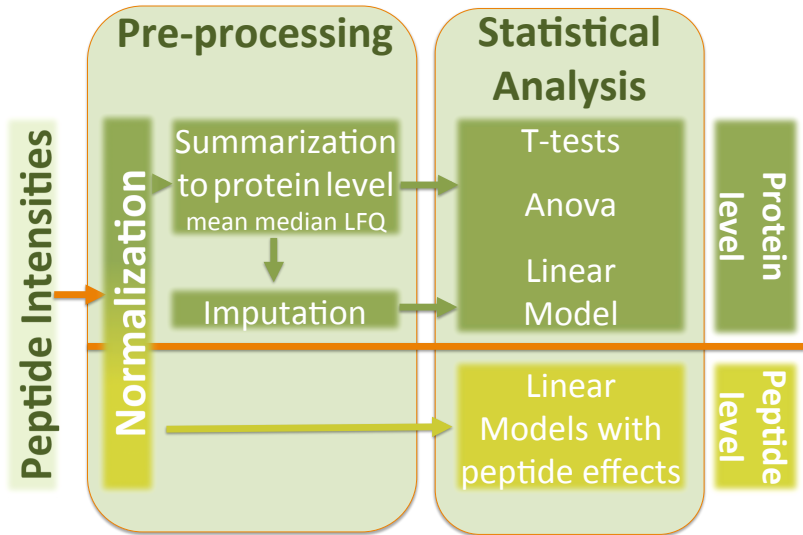
We need information on protein level!



Label-free Quantitative Proteomics Data Analysis Pipelines



Label-free Quantitative Proteomics Data Analysis Pipelines



CPTAC Spike-in Study

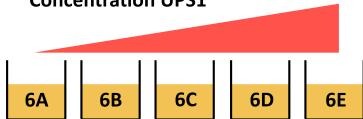
Digested
UPS1 protein mix



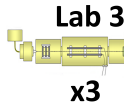
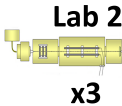
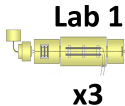
Digested
yeast proteins



Concentration UPS1



5 spike-in concentrations: 6A to 6E



- Same trypsin-digested yeast proteome background in each sample
 - Trypsin-digested Sigma UPS1 standard: 48 different human proteins spiked in at 5 different concentrations (treatment A-E)
 - Samples repeatedly run on different instruments in different labs
 - After MaxQuant search with match between runs option
 - 41% of all proteins are quantified in all samples
 - 6.6% of all peptides are quantified in all samples
- **vast amount of missingness**

Preprocessing

- Typical preprocessing steps
 - ① Filtering
 - ② Log-transformation
 - ③ Normalization
 - ④ (Summarization)

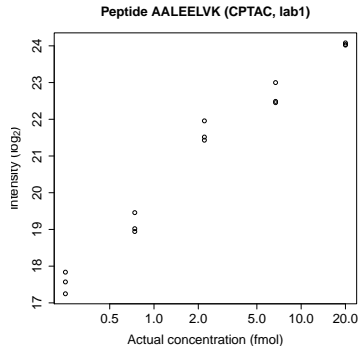
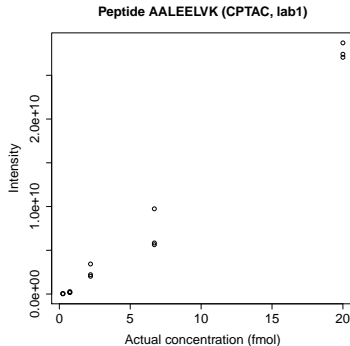
- Many methods exist

Filtering

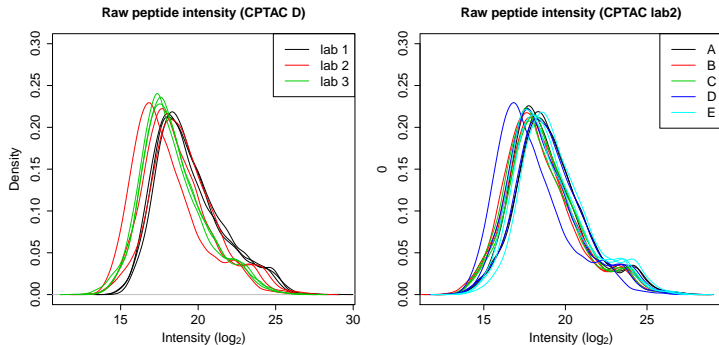
- Reverse sequences
- Only identified by modification site (only modified peptides detected)
- Razor peptides: non-unique peptides assigned to the protein group with the most other peptides
- Contaminants
- Peptides few identifications
- Proteins that are only identified with one or a few peptides

- Filtering does not induce bias if the criterion is independent from the downstream data analysis!

Log-transformation



Variability more equal upon log transformation: often multiplicative error structure of intensity-based read-outs



Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

- Considerable effects between and within labs for replicate samples
 - Considerable effects between samples with different spike-in concentration
- Normalization is needed

Mean or median?

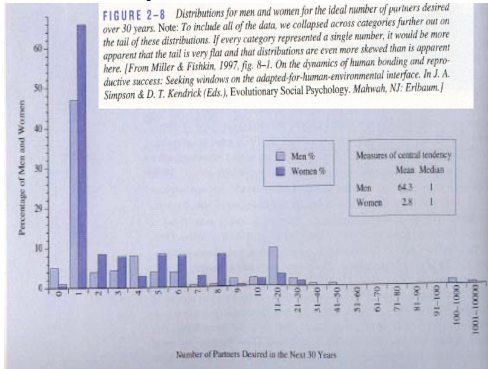
- Over a period of 30 years males desire to have on average 64.3 partners and females 2.8. (Miller and Fishkin, 1997)

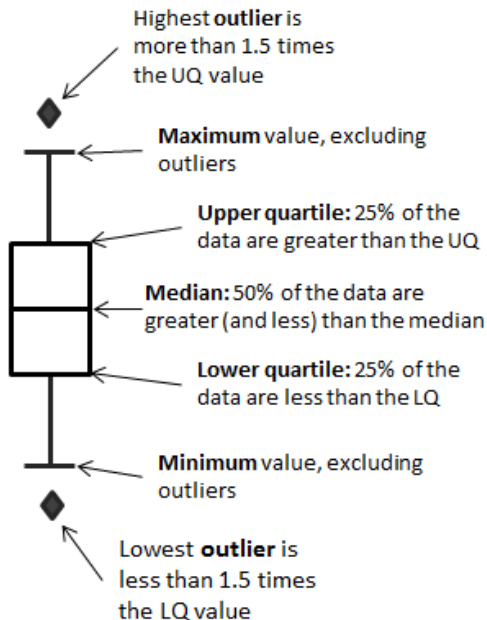
Mean or median?

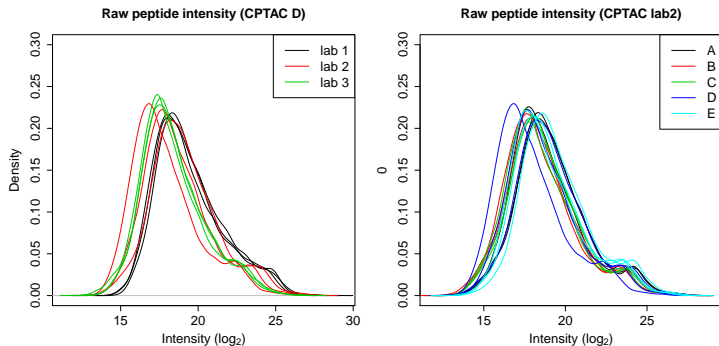
- Over a period of 30 years males desire to have on average 64.3 partners and females 2.8. (Miller and Fishkin, 1997)
- Over a period of 30 years males, is the median of the number of desired partners is 1 for both males and females. (Miller and Fishkin, 1997)

Mean or median?

Mean is very sensitive to outliers!

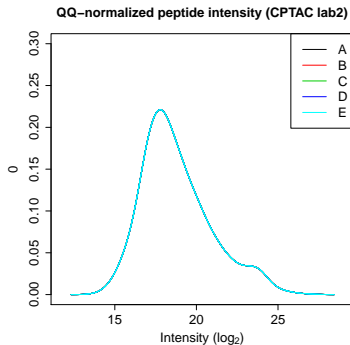
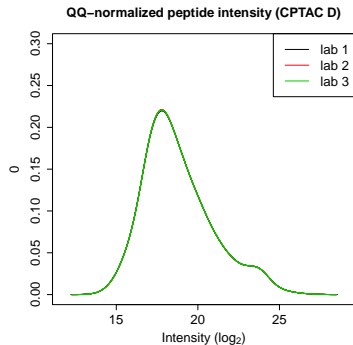






Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

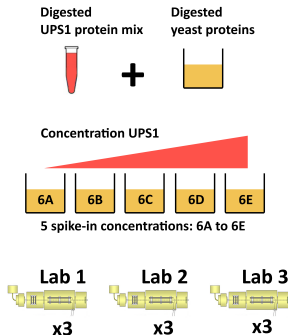
- Considerable effects between and within labs for replicate samples
 - Considerable effects between samples with different spike-in concentration
- Normalization is needed



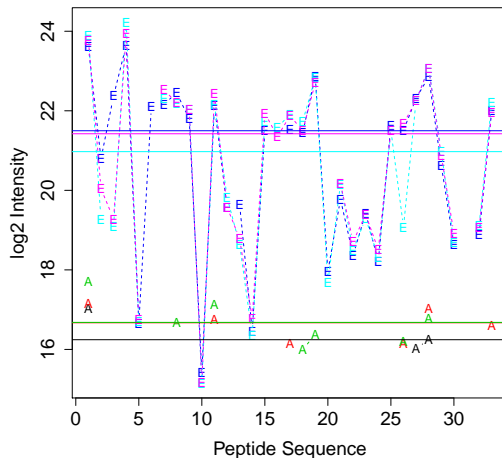
Even in very clean synthetic dataset (same background, only 48 UPS proteins can be different) the marginal peptide intensity distribution across samples can be quite distinct

- Considerable effects between and within labs for replicate samples
 - Considerable effects between samples with different spike-in concentration
- Normalization is needed, e.g. **quantile normalization**

Summarization

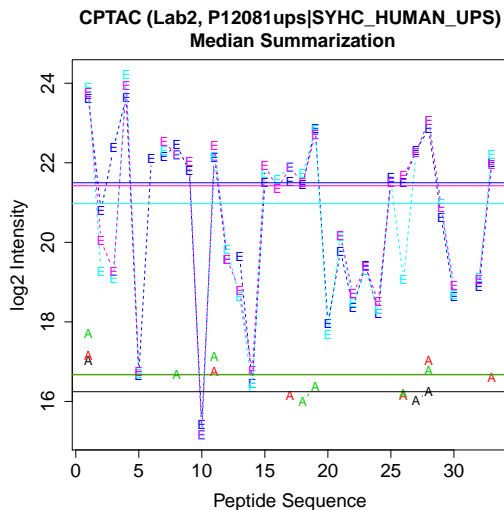


CPTAC (Lab2, P12081ups|SYHC_HUMAN_UPS)
Median Summarization



Summarization

- Strong peptide effect
- Unbalanced peptide identification
- Summarization bias
- Different precision of protein level summaries



MaxLFQ summarization

a

>P63208

MPSIKLQSSDGEIFEVDVEIAKQSVTIKTMLEDLGMDDEGDD
 DPVPLPNVNAAILKKVIQWCTHHKDDPPPPEDDENKEKRTDD
IPVWDQEFLEKVDQGTFLFELILAAANYLDIKGLLDVTCKTVANM
IKGKTPEEIRKTFNIKNDFTEEEEAQVRKENQWCEEK

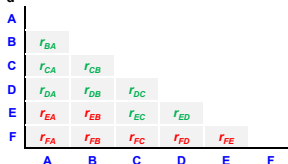
b

Peptide species	Sequence	Charge	Mod.
P ₁	LQSSDGEIFEVDVEIAK	2	-
P ₂	LQSSDGEIFEVDVEIAK	3	-
P ₃	RTDDIPVWDQEFLEK	2	-
P ₄	TVANMIK	2	-
P ₅	TVANMIK	2	Oxid.
P ₆	TPEEIRK	3	-
P ₇	NDFTEEEEAQVR	2	-

c

Sample	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆	P ₇
A		+				+	
B		+	+			+	
C	+	+	+	+		+	+
D	+	+		+		+	+
E		+		+			+
F		+			+		

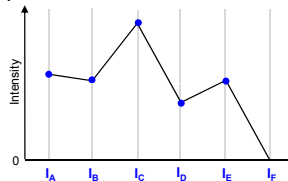
d



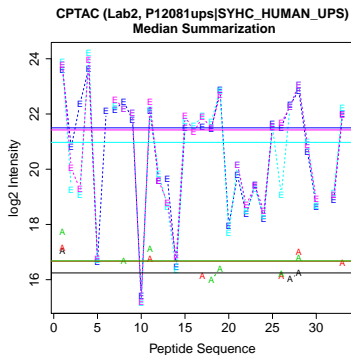
e

$r_{BA} = I_B / I_A$	$r_{CA} = I_C / I_A$	$r_{CB} = I_C / I_B$
$r_{DA} = I_D / I_A$	$r_{DB} = I_D / I_B$	$r_{DC} = I_D / I_C$
$r_{EC} = I_E / I_C$	$r_{ED} = I_E / I_D$	$I_F = 0$

f

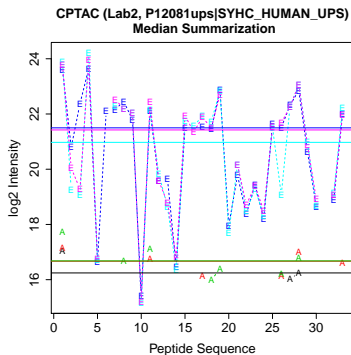


Peptide based model



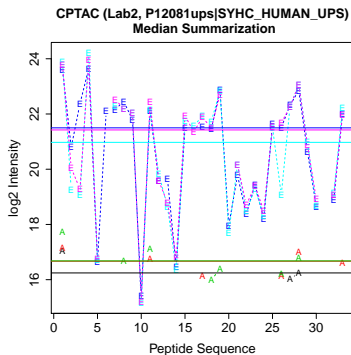
- ① y_{ip} : log₂ intensity for peptide p of a particular protein in sample i

Peptide based model



- 1 y_{ip} : log2 intensity for peptide p of a particular protein in sample i
- 2 Protein by protein analysis of peptide data with linear model

Peptide based model

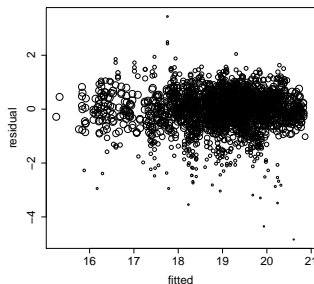
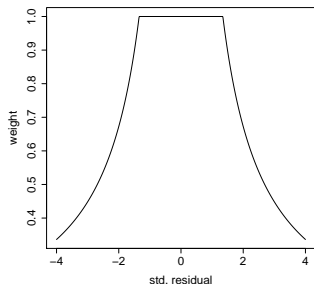


- 1 y_{ip} : log2 intensity for peptide p of a particular protein in sample i
- 2 Protein by protein analysis of peptide data with linear model

$$\begin{array}{c}
 \text{peptide level} \\
 y_{ip} = \beta_p^{\text{pep}} + \epsilon_{ip}
 \end{array}
 +
 \begin{array}{c}
 \text{protein level} \\
 \beta_i^{\text{sample}}
 \end{array}$$

Robust estimation using observation weights

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



- Iteratively fit model with observation weights $w(\epsilon_{ip})$

$$\operatorname{argmin}_{\beta_{1\dots P}^{\text{pep}}, \beta_{1\dots n}^{\text{samp}}} \left[\sum_{i=1}^n \sum_p^P w(\epsilon_{ip}) (y_{ip} - \beta_p^{\text{pep}} - \beta_i^{\text{samp}})^2 \right]$$