



Quantitative Proteomics

A short introduction

Martin Wells

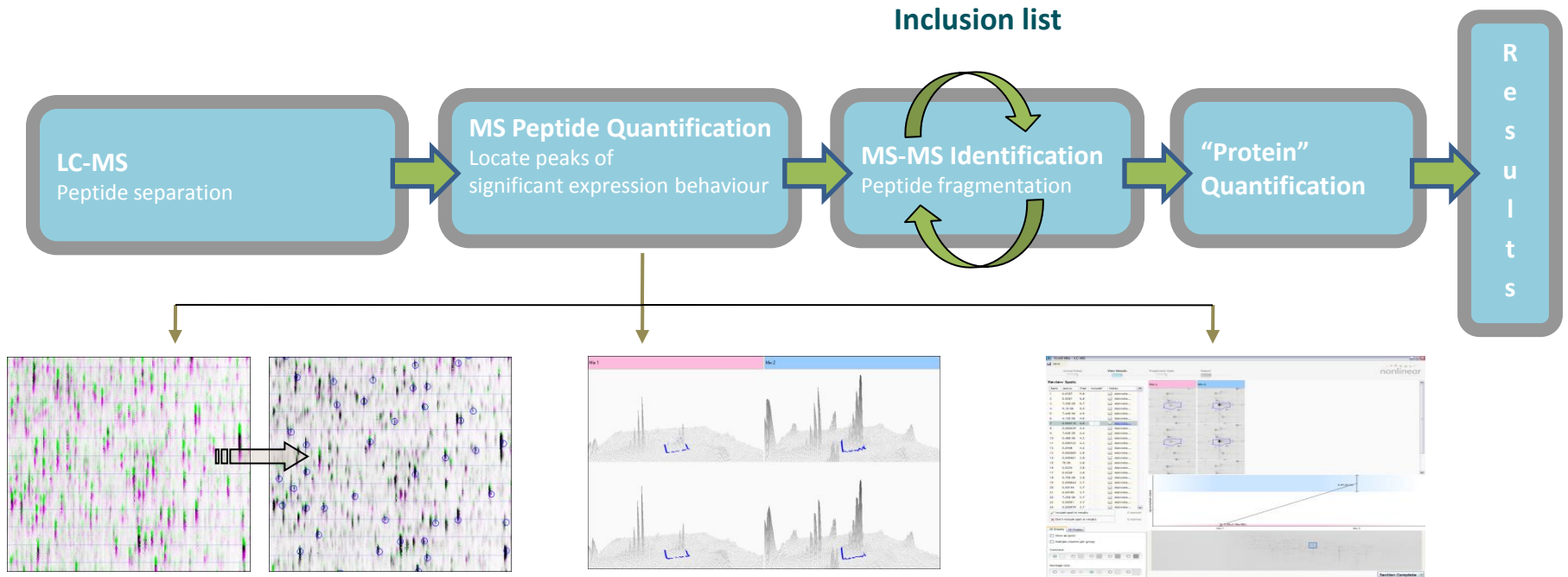
15th – 16th July 2010
Hinxton, Cambridge

Quantitative Proteomics

- Gel or gel free
- Labelled or label free
- Relative or absolute

**Gel free, label free, relative
quantification**

Analytical workflow



LC-MS data alignment algorithms corrects for the positional bias introduced by the LC

Co-detection and relative quantification of all peptide ions and statistical tools to define subsets of experimental interest.

Integration with search engines to identify peptide ions. Create inclusion lists for repeat analysis of peptides unidentified.
Combining of all data at the protein level and calculation of the protein quantification.

Major contributors to data variance

- Technical noise
 - e.g. experimentally introduced bias
- Analytical inadequacies
 - e.g. incorrect identification and quantification of experiment data
- Biological variation
 - how do we handle this?



These obstacles result in a loss of power and, therefore, a reduction in our ability to discover significant expression changes

What is the goal of differential expression analysis?

“To identify the best areas that warrant further investigation as rapidly, objectively and reliably as possible.”

We need to:

1. Process large data files efficiently enabling the confident selections of features for further investigation
2. Limit user subjective handling and interpretation of the data to enable reproducible analysis and results
3. Increase confidence in the resulting protein lists through improved data quality. Missing data is the primary limitation to the application of multivariate statistical techniques so addressing this would significantly impact data reliability and confidence.

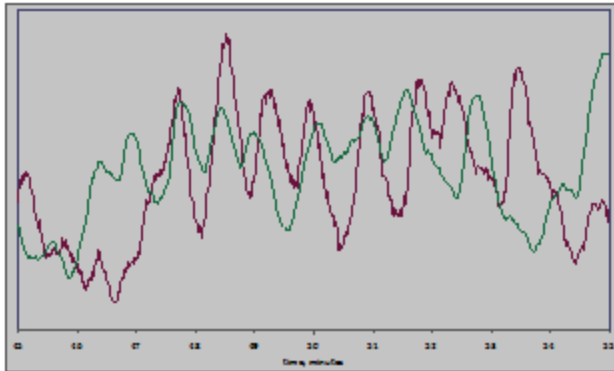
Data alignment

Assumptions

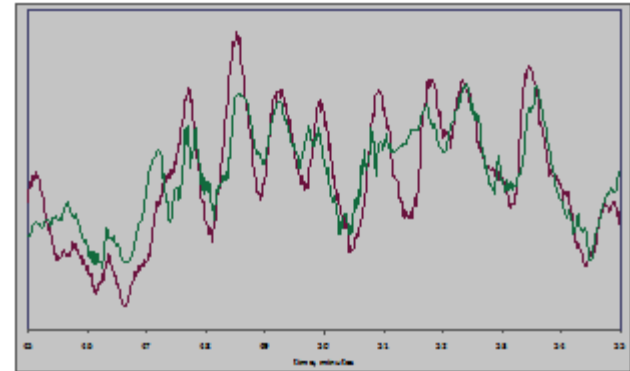
- A peptide ion, if present, should theoretically always appear in the same position under identical experimental conditions
- Most differential expression experiments have a high degree of common peptide ions between groups

As each peptide ion has a unique position can we, using the common peptide ions, realign the data so that we undo the positional bias that has been introduced by the chromatography?

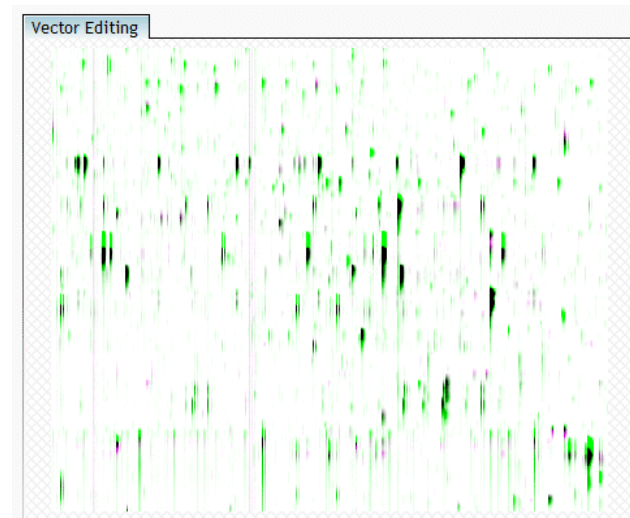
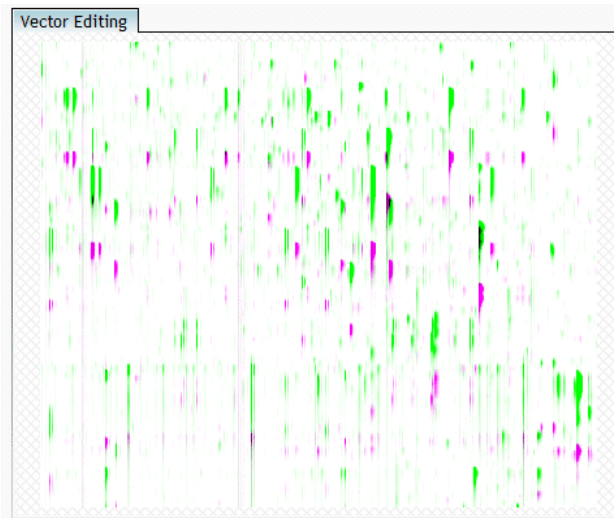
Data alignment



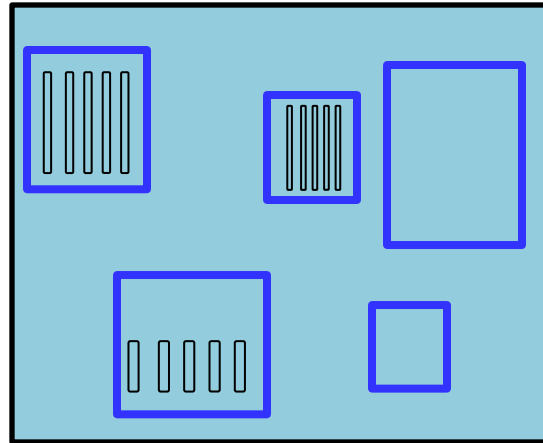
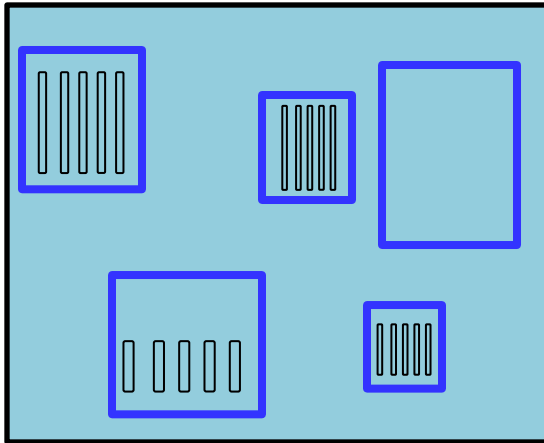
Section prior to alignment



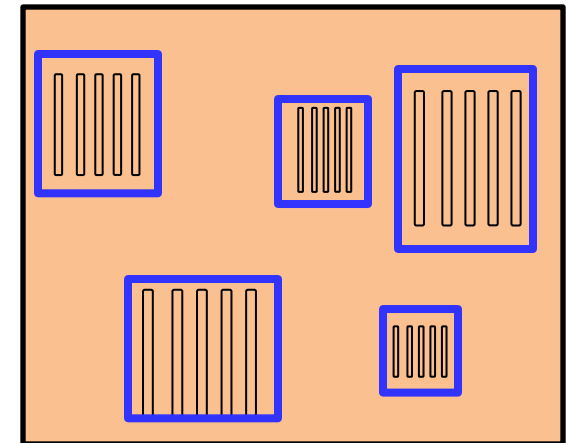
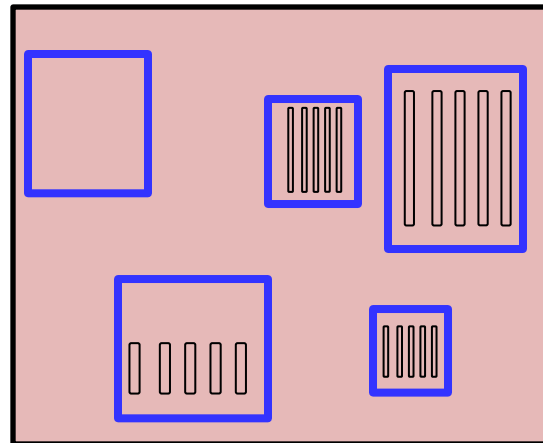
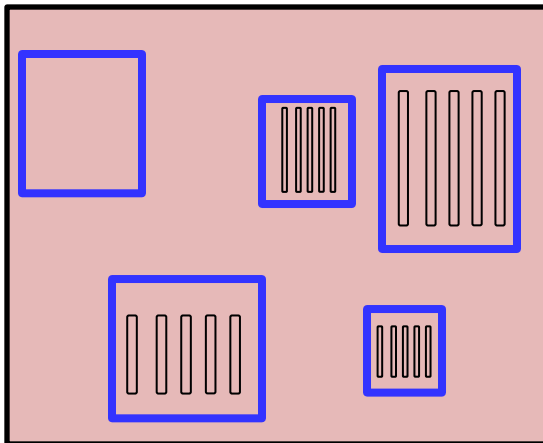
Section after alignment



Feature detection – peak picking

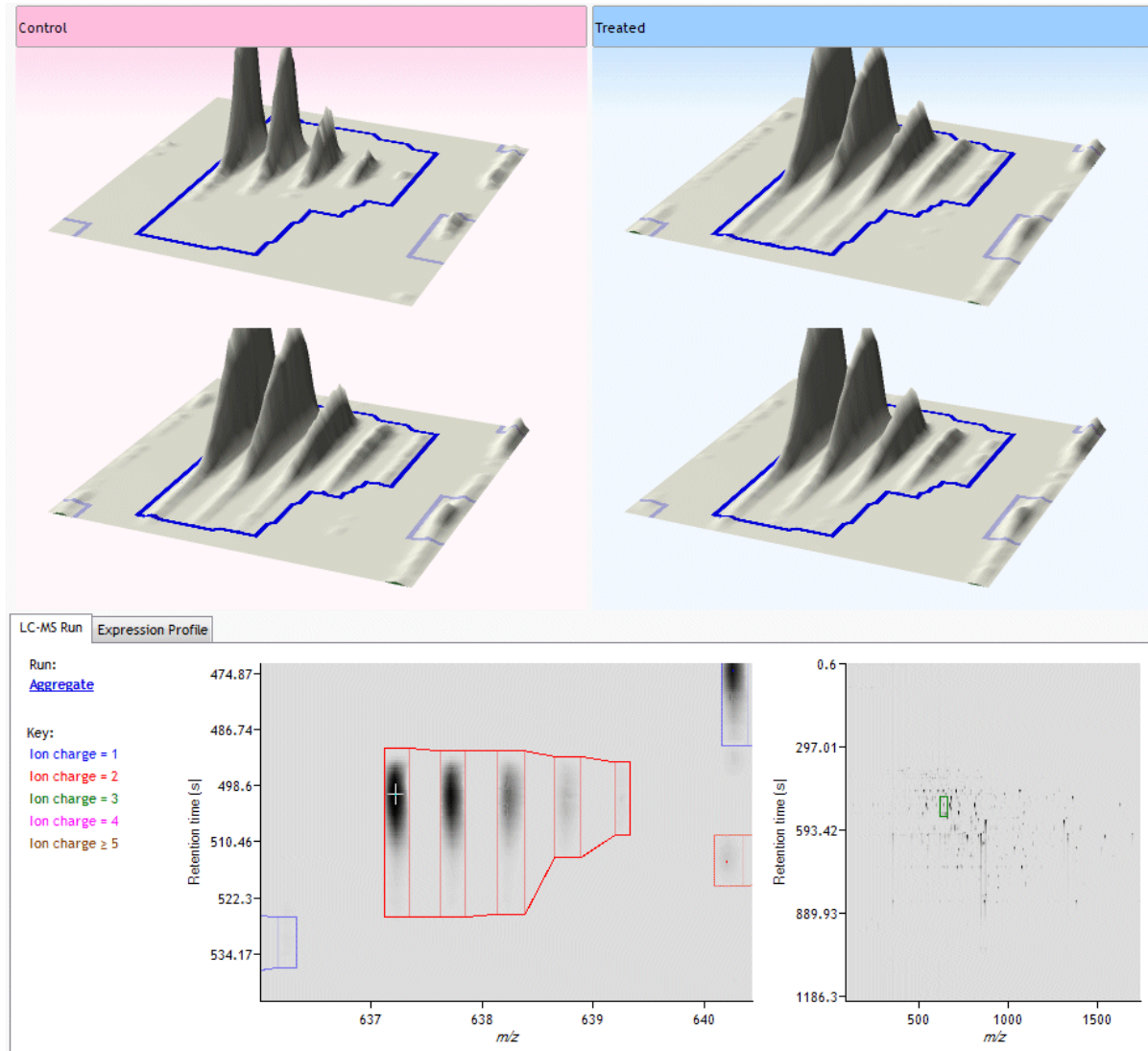


Mapping the detection to all runs avoiding missing data



Aggregate co-detection

Peptide ion raw abundance



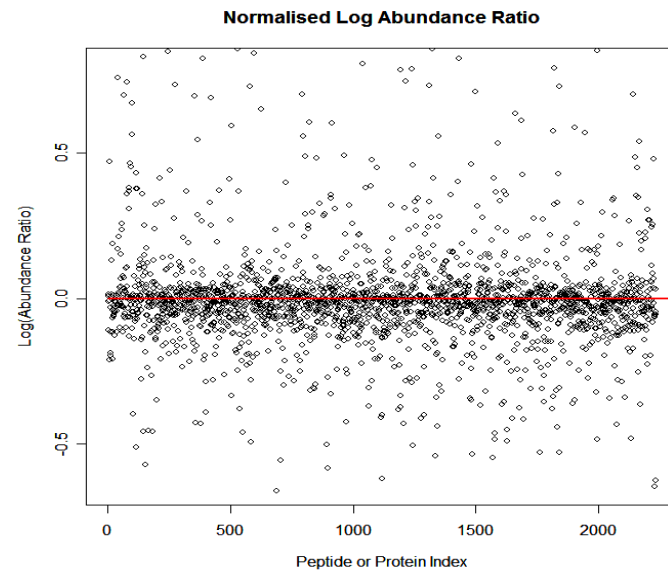
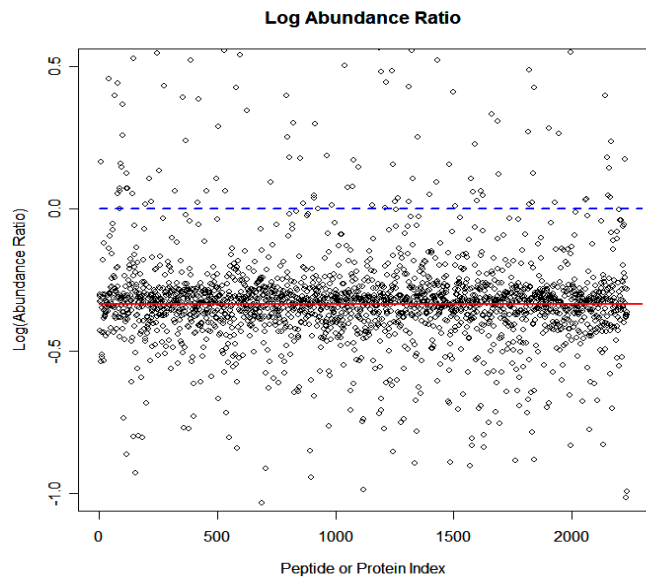
Data normalisation

Normalisation is required to calibrate between different sample runs.

Total ion current has been assumed to be equal across all samples. However without missing data an alternative approach is possible.

Calculating log abundance ratios for every feature gives a more robust assessment of the gain factor.

Of course, the base assumption still remains that enough features should NOT be changing in abundance.

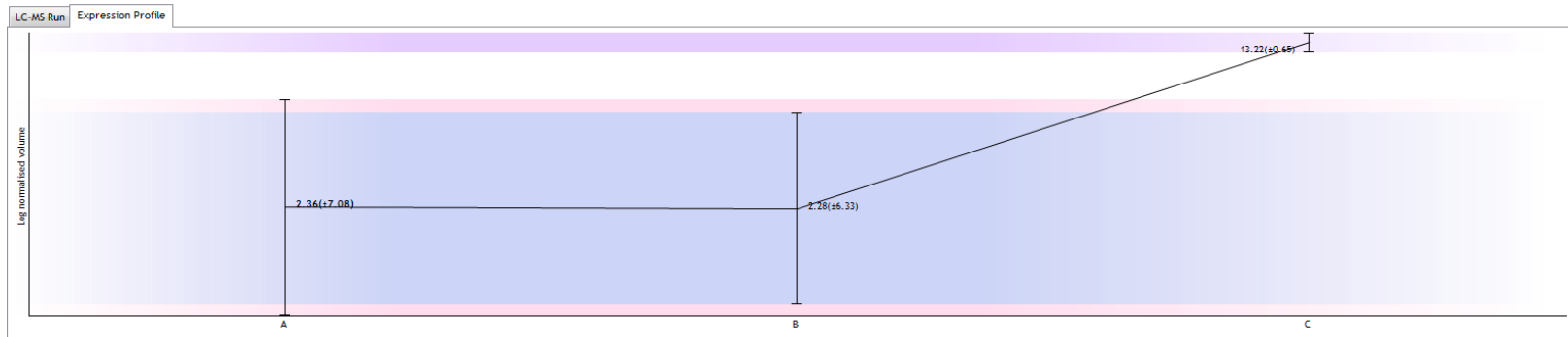


Relative normalised abundances.



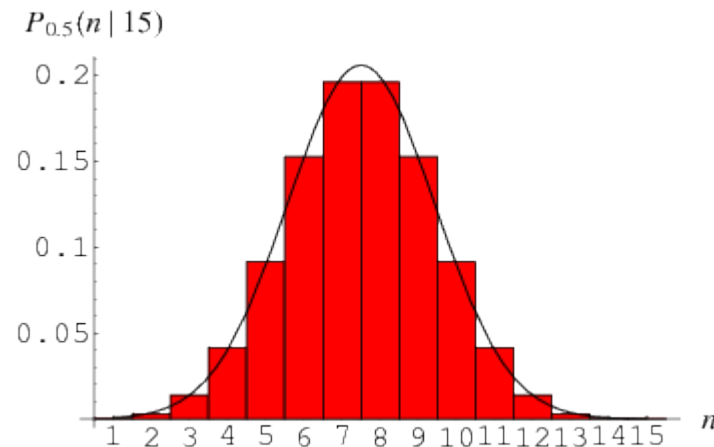
3 groups with 3 replicates per group.

Graphical expression profile showing group means with error bars representing 95% confidence bounds.



Data transformation

The data must meet certain prerequisites for the statistical test applied.
ie The data is normally distributed and of equal variance.



To achieve this a **data transformation** is required.

An ArcSinh transformation of the normalised data is most suitable particularly at very low values when compared to other transformations eg log.

Therefore all statistical analysis is conducted using the ArcSinh normalised volumes

Uni and multivariate statistics

The complete data table of quantitative information, without missing data, enables the robust and confident application of multivariate statistics to support traditional univariate tests.

- p-values (typically $p < 0.05$ – ie 5% type I error)
- PCA – Quality control overview of the variances observed in your data.
- Power analysis per feature to reduce type II errors
- Predictive power for predicting how many samples are needed to be confident you are not missing something
- False discovery rate correction – Q values to reduce type 1 errors
- Correlation analysis of peptide ions
See which peptides are highly correlated within and across groups.

Statistical Power

In addition to the standard p-value filtering of data, by addressing the issue of missing data, we are now able to accurately calculate the statistical power for each feature.

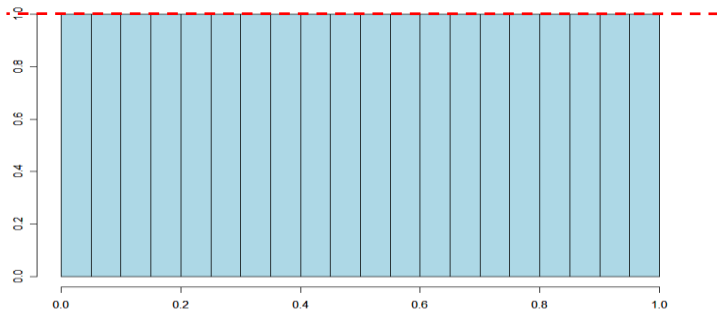
Power – the error of failing to observe a difference when in truth there is one, thus indicating a test of *poor sensitivity*

This enables secondary filtering of the data to define a subset of features that have both low p-value (eg <0.05) and high power (eg $>80\%$).

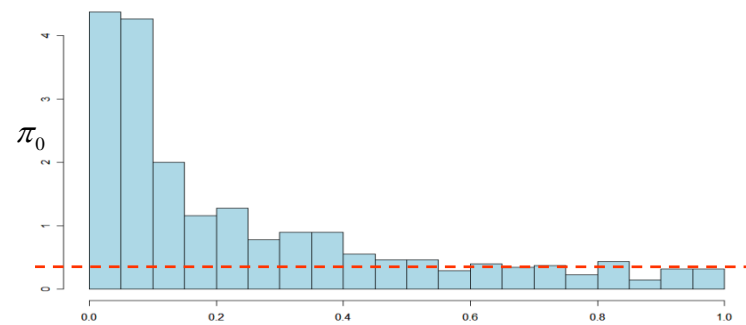
You can calculate that you have run sufficient samples to observe a difference if one really exists or reduce your “discoveries” to features where confidence is high. (eg $>80\%$)

False Discoveries and Q-values

- False discoveries occur when we are observing a difference when in truth there is none, thus indicating a test of *poor specificity*. *Q-values* are a tool to reduce such *false discoveries*
- False Discovery Rate (FDR) assumes a uniform distribution of p-values but ...



Theoretical p-value distribution



Actual p-value distribution

- q-values take the into account the actual p-value distribution factor

$$\tilde{q}_m = p_m$$
$$\tilde{q}_i = \min(\tilde{q}_{i+1}, \pi_0 * (m/i) * p_i), i = m-1, \dots, 1$$

[1] Storey, J.D. and Tibshirani, R. (2003) Statistical significance for genomewide studies. Proc. Natl Acad. Sci. USA, **100**, 9440–9445.

[2] Hedenfalk, I., Duggan, D., Chen, Y. D., Radmacher, M., Bittner, M., Simon, R., Meltzer, P., Gusterson, B., Esteller, M., Kallioniemi, O. P., et al. (2001) N. Engl. J. Med. **344**, 539–548.

Database result integration

MSMS fragment scans are extracted and compiled into file formats compatible with most search engines.

- Mascot
- Sequest
- Phenyx
- PLGS
- Peaks Studio
- Proteome Discoverer

The qualitative results from can then be imported back into the software for integration with the quantitative data.

Data is then automatically collated at the protein level and the protein quantification is then calculated based on the combined impact of all peptide ion measurements.

Tools are available to remove non-unique peptides from the data and highlight peptides not in consensus with other peptide abundances. Eg PTM's

Protein results

LC-MS Tutorial with ID's - LC-MS

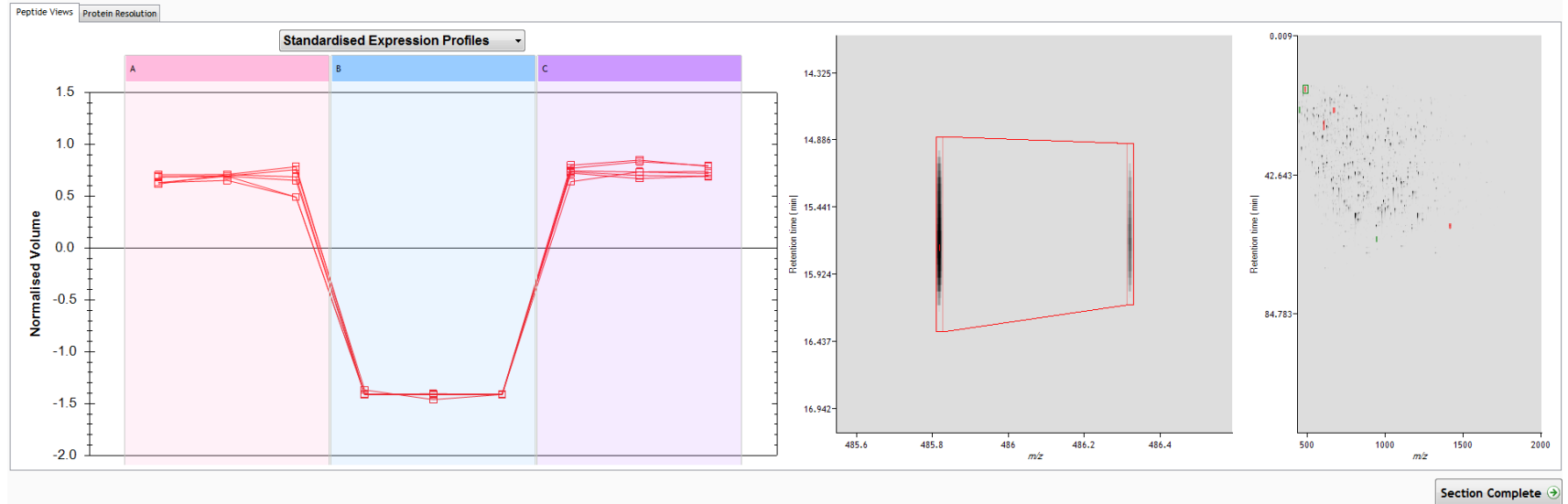
File

LC-MS Data Import Reference Run Selection Alignment Filtering Group Setup View Results Progenesis Stats Protein Search Protein View Report

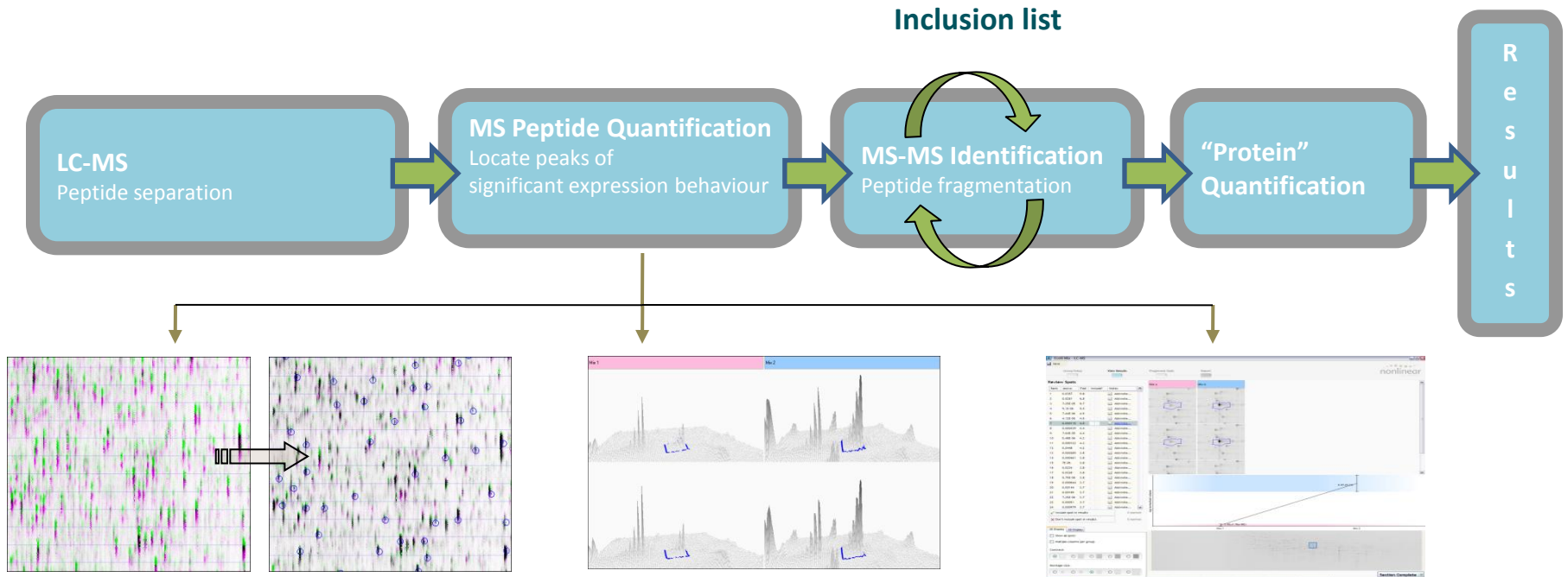
nonlinear DYNAMICS

Protein: **gi|196163602 cell surface protein (S-layer precursor protein) [Clostridium difficile QCD-23m63]**

Accession	Peptides	Mass	p-value*	Fold	Abundance	Score	#	Score	Hits	Mass	RT (mins)	Charge	Tags	Abundance	Proteins	Peptide Sequence
gi 170632806	3	77784	5.51E-10	1.06E+03	3.51E+07	247.80	5	78.42	48	1206.66	26.2	2	✓	2.21E+07	6	VYLAGGVNSISK
gi 170632792	3	76101	7.89E-10	798	2.32E+07	195.94	376	26.63	16	2822.35	58.1	2	✓	3.81E+06	3	AKTLSSDASDFLGNAQVDIIIGGENSVSK
gi 196163602	4	74575	2.96E-10	899	2.61E+07	262.52	576	5.9	1	2822.50	62.8	3	✓	1.16E+06	3	AKTLSSDASDFLGNAQVDIIIGGENSVSK
gi 188587165	1	28390	9.27E-11	Infinity	1.71E+05	82.60	2033	47.88	10	1335.74	22.5	3	✓	8.77E+04	5	KVYLAGGVNSISK
gi 168713980	27	79874	1.74E-10	1.51E+03	3.18E+08	1936.78	2131	55.32	20	1335.74	22.4	2	✓	1.47E+05	5	KVYLAGGVNSISK
gi 163762164	2	164000	2.32E-10	Infinity	1.37E+05	46.92	8037	48.37	9	969.62	15.8	2	✓	7.59E+03	3	KAPLLLTSK
gi 146295117	2	30682	9.32E-11	Infinity	3.26E+05	45.76										
gi 145953274	27	80379	1.74E-10	1.51E+03	3.18E+08	1886.53										
gi 145953272	10	66427	1.86E-12	Infinity	3.84E+06	846.81										
gi 145953270	5	66927	8E-13	3.24E+07	3.78E+05	283.79										
gi 145953209	2	51355	8.63E-10	764	2.22E+07	45.80										
gi 126700400	1	72997	1.38E-08	Infinity	1.88E+04	81.59										
gi 126698643	1	40834	7.32E-11	Infinity	1.27E+05	106.05										
gi 125718796	2	39307	3.23E-13	Infinity	1.74E+05	46.31										



Analytical workflow



LC-MS data alignment algorithms corrects for the positional bias introduced by the LC

Co-detection and relative quantification of all peptide ions and statistical tools to define subsets of experimental interest.

Integration with search engines to identify peptide ions. Create inclusion lists for repeat analysis of peptides unidentified. Combining of all data at the protein level and calculation of the protein quantification.

Practical Session

Experiment

- Membrane proteins extracted from 2 different strains of Clostridium Difficile.
- Three separate cultures were prepared for each strain and analysed on the LC-MS. (Dionex Ultimate 3000 and Thermo Orbitrap XL)
- Analyse the 6 resulting files, comparing Strain “A” to strain “C”
- A database search has been performed using mascot. Import the identification results, removing hits with a mascot scores <30 and hits not containing “difficile” in the protein description.
- After removing any non-unique peptides from the quantification, report the fold change and p-values of the proteins with >4 unique peptides.
- For the purpose of this exercise hypothetical proteins can be excluded from the results.

Practical Session

Expected results...

6 Proteins all up regulated in strain C (excluding hypothetical proteins).

Protein Description	p-value	Fold change	No. Unique peptides
ABC transporter, substrate Binding lipoprotein	1.3E-05	23	12
Cell surface protein	0.0002	3	5
Enolase	3.7E-05	3.5	5
Putative 5-nitroimidazole Reductase	3.7E-06	31	5
Glyceraldehyde-3-phosphate Dehydrogenase	2.3E-06	8	4
Cell wall protein V	3.4E-07	100	4