

Using R and Bioconductor for proteomics data analysis

L. Gatto^{*1}, L.M. Breckels¹, S. Gibb², A. Christoforou¹ and K. S. Lilley¹

¹Cambridge Centre for Proteomics, Department of Biochemistry, University of Cambridge, UK

²Institut für Medizinische Informatik, Statistik und Epidemiologie, Universität Leipzig, Germany

^{*}lg390@cam.ac.uk – <http://www.bio.cam.ac.uk/proteomics/>

Introduction

R and Bioconductor have had a tremendous impact on the quality of genomics data analysis [1], demonstrating that extraction of relevant and biologically meaningful information from high-throughput data, required investing time and effort in the exploration and analysis of the data.

Even well-known and respected leader in proteomics agree that it lies 10 years behind genomics. There are several valid reasons for this, including the chemical complexity of proteins, the technical complexity of the instrumentation (in particular mass-spectrometry - MS) and the vast possibilities in the study of proteins. An often overseen albeit essential component of this failure is arguably the software that is promoted inside the proteomics community. Computational proteomics researcher who value quality software, comprehensive data analysis and reproducible research ought to illustrate how more flexible and advanced tools can effectively be used and demonstrate their advantages. Here, we illustrate some examples of proteomics data analysis in R .

Working with raw data

The proteomics community has developed a range of data standards and formats for MS data (the latest being mzML) to overcome the shortcomings or closed, binary vendor-specific formats. One of the main projects that implement parsers for the XML-based open formats is the C++ proteowizard project [2], which is interfaced by the mzR Bioconductor package using the Rcpp package.

```
library("mzR")
fname <- dir(system.file(package = "MSnbase", dir = "extdata"),
             full.name = TRUE, pattern = "mzXML$")
ms <- openMSfile(fname)
```

The resulting ms object is a file handle that allows fast direct and random access to the individual spectra. mzR is used by a variety of other packages like xcms, MSnbase, RMassBank and TargetSearch.

Challenges Further improve support of raw MS data and develop the range of supported formats, in particular identification (mzIdentML) and quantitation (mzQuantML) formats.

Labelled quantitation

The same raw data file can be imported in a convenient higher level container and directly processed, plotted, quantified and normalised with the MSnbase software.

```
exp <- readMSData(fname, verbose = FALSE)
set <- quantify(exp, method = "trap", reporters = iTRAQ4,
               verbose = FALSE, parallel = TRUE)
head(exprs(set), n = 3)
```

:	iTRAQ4.114	iTRAQ4.115	iTRAQ4.116	iTRAQ4.117
: X1.1	4483	4874	6743	4601
: X2.1	1918	1418	1118	1582
: X3.1	15211	15296	15593	16551

```
plot(exp[["X3.1"]], full = TRUE, reporters = iTRAQ4)
```

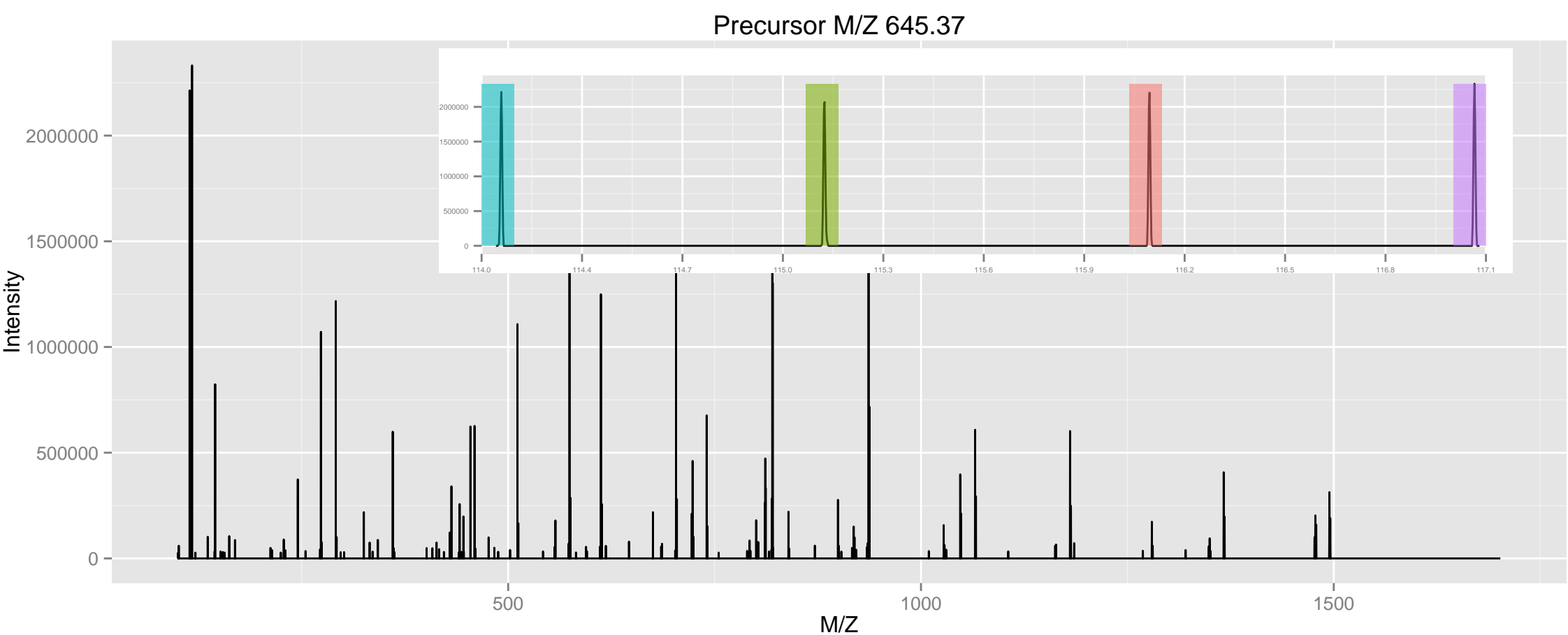


Figure 1: MS² spectrum of an iTRAQ 4-plex experiment showing the 4 isobaric reporter ions.

Challenges Although labelled MS² quantitation is well supported with MSnbase and isobar, metabolic labelling techniques like N¹⁵ or SILAC still need to be developed.

Label-free quantitation

MS^e data independent acquisition

Peptide identification

Conclusions and perspectives

The felxibility of the R environment and the breath of available packages is sometimes daunting for newcomers and introductory points of entry are wel-come. The RforProteomics package [3] [<https://github.com/lgatto/RforProteomics>] ought to assume such roles. For this, RforProteomics should be a collaborative project and contribution through the github repository are encouraged.

Despite well known advantages in terms of statistical analyses of data and some unique software for proteomics and mass-spectrometry data analysis, there remains a lot of efforts and work to be done for R/Bioconductor to become a complete framework for proteomics data processing. These efforts should be tackled by a group of developers. It is our hope that the RforProteomics will be a helpful targeted introduction to new users and motivate collaborative development of package developers.

¹ Gentleman et al. Bioconductor: open software development for computational biology and bioinformatics. Genome Biol. 2004 PMID: 15461798.

² Chambers et al. A cross-platform toolkit for mass spectrometry and proteomics. Nat Biotechnol. 2012 PMID: 23051804.

³ Gatto L, Christoforou A. Using R and Bioconductor for proteomics data analysis. Biochim Biophys Acta. 2013 PMID: 23692960.