

Using R and Bioconductor for proteomics data analysis

Laurent Gatto^{1,*}, Andy Christoforou

*Cambridge Centre for Proteomics, Department of Biochemistry
University of Cambridge, Tennis Court Road, Cambridge, CB2 1QR, UK*

Abstract

This review presents how R, the popular statistical environment and programming language, can be used in the frame of proteomics data analysis. A short introduction to R is given, with special emphasis on some of the features that make R and its add-on packages a premium software for sound and reproducible data analysis. The reader is also advised on how to find relevant R software for proteomics. Several use cases are then presented, illustrating data input/output, quality control, quantitative proteomics and data analysis. Detailed code and additional links to extensive documentation are available in the freely available companion package RforProteomics.

Keywords: software, mass spectrometry, quantitative proteomics, data analysis, statistics, quality control

1. Introduction

Proteomics is evolving at a rapid pace [1] and updates in technologies and instruments applied to the study of bio-molecules, such as proteins or metabolites, require proper computational infrastructure [2]. A broad diversity of complementary tools for data processing, management, visualisation and analysis have already been offered to the community and reviewed elsewhere [3, 4]. The work presented here focuses on a particular type of software, namely R [5], and the add-on *packages* that enable extension in its functionality and scope, and their usefulness to the analysis of proteomics data.

*Corresponding author

Email addresses: lg390@cam.ac.uk (Laurent Gatto), ac587@cam.ac.uk (Andy Christoforou)

¹Tel: +44 1223 760253 Fax: +44 1223 333345

9 R is an open source statistical programming language and environment, originally created
10 by Ross Ihaka and Robert Gentleman [6] at the University of Auckland and, since the mid-
11 1997, developed and maintained by the R-core group. Originally utilised in an academic
12 environment for statistical analysis, it is now widely used in public and private sector in a
13 broad range of fields [7], including computational biology and bioinformatics. The success
14 of R can be attributed to several features including flexibility, a substantial collection of
15 good statistical algorithms and high-quality numerical routines, the ability to easily model
16 and handle data, numerous documentation, cross-platform compatibility, a well designed
17 extension system and excellent visualisation capabilities to list some of the more obvious
18 ones [8]. These are some of the requirements that need to be fulfilled to tackle the complexity
19 and high-dimensionality of modern biology.

20 The focus of R itself is and remains centred around statistics and data analysis. Function-
21 ality can however be extended through third-party packages, which bundle a coherent set of
22 functions, documentation and data to address a specific problem and/or data type of inter-
23 est. The Bioconductor project² [9], initiated by Robert Gentleman, has a specific focus on
24 computational biology and bioinformatics and represents a central repository for hundreds
25 of software, data and annotation packages dedicated to the analysis and comprehension of
26 high-throughput biological data, and promoting open source, coordinated, cooperative and
27 open development of interoperable tools. The development and distribution of new packages
28 is a very dynamic and important aspect of the R software itself. Adherence to good devel-
29 opment practice is crucial and enforced by the R package development pipeline through a
30 built-in checking mechanism, ensuring, among other things, proper package installation and
31 loading, package structure, code validity and correct documentation. In addition, package
32 development also provides multiple opportunities for unit and integration testing as well as
33 reproducible research [10, 11, 12, 13, 14] through the mechanism of literate programming
34 [15] and Sweave [16] or knitr [17] vignettes, which is crucially important from a scientific
35 perspective.

²<http://bioconductor.org/>

36 Packages can be submitted to the main central repository, the Comprehensive R Archive
37 Network (CRAN) or to Bioconductor, which provides its own repository, to assure tighter
38 software interoperation. In addition, any developer can easily set up private or public
39 CRAN-style systems. Software management can become a tedious task when thousands of
40 packages are distributed, many of which depend on each other and interoperate in complete
41 pipelines. In R, this has been solved by providing dedicated package repositories as well as
42 straightforward installation and updating mechanisms.

43 Most importantly, R and many packages are regarded as quality software [18]. They are
44 aimed at users who want to explore and comprehend complex data for which there is often
45 no predefined recipe. It is also a research tool to tackle new questions in innovative ways.
46 The Bioconductor project, for example, has had a substantial impact on the field of microar-
47 rays through multi-disciplinary and cooperative method development and implementation,
48 paving best practises for the current development of state-of-the-art high throughput ge-
49 nomics data analysis and comprehension. With respect to R's contribution to other areas
50 of bioinformatics and computational biology, it has also a lot to offer to proteomics. Biolo-
51 gists and proteomicists can gain immensely from autonomous data exploration and analysis.
52 Bioinformaticians working in computational proteomics can use R and specialised packages
53 as an independent analysis and research framework or employ them to complement existing
54 pipelines.

55 This manuscript presents a brief overview of some applications of the R software to
56 the analysis of MS-based quantitative proteomics data. We will review compliance of R
57 with open proteomics data standards, input/output capabilities, quantitation pipelines for
58 label-free and labelled quantitation, quality control, quantitative data analysis and relevant
59 annotation infrastructure. The review is accompanied by a package, `RforProteomics`, that
60 provides the code to install a selection of relevant tools to reproduce and adapt the examples
61 described below. Installation instruction are provided on the package's web page³. Once
62 installed, the package is loaded with the `library` function as shown below, to make its

³<http://lgatto.github.com/RforProteomics/>

63 functionality available.

```
> library("RforProteomics")  
  
This is the 'RforProteomics' version 1.0.1.  
Run 'RforProteomics()' in R or visit  
'http://lgatto.github.com/RforProteomics/' to get started.
```

64 2. Using R in proteomics

65 2.1. Finding relevant software

66 R is a very dynamic *ecosystem* [19, 20] – yearly R and bi-annual Bioconductor releases,
67 exponentially growing number of available packages [21], numerous active mailing lists and a
68 community of hundreds of thousands of active users and developers in private and corporate
69 environment [7]. There are currently thousands of packages available through the official
70 repositories, and new packages are published, discontinued or replaced by new, more elabo-
71 rate alternatives on a daily basis. Providing an up-to-date and exhaustive list of packages
72 is unachievable, even for a specified area of interest like proteomics, and would undoubtedly
73 be out-dated too quickly to be useful. Dedicated pages are available however, that allow one
74 to obtain an overview of some of the available packages in a specific area. CRAN maintains
75 topic task views⁴, which are curated and maintained by experts. Each view provides a sum-
76 mary and some guidance on some of the growing number of CRAN packages that are useful
77 for a certain topic. As of this writing, the Chemometrics and Computational Physics view
78 features a total of 67 packages, some of which are dedicated to mass spectrometry and will
79 be described later. The Bioconductor project provides a set of dedicated keywords to cate-
80 gorise packages, called *biocViews*, that can be explored interactively⁵. For proteomics, most
81 relevant candidates are *MassSpectrometry* (in the *Software/AssayTechnology* view with 21
82 packages) and *Proteomics* (in the *Software/BiologicalDomain* view, 35 packages), although
83 numerous data analysis and annotation packages in other categories provide invaluable sup-
84 port, some of which will also be demonstrated below.

⁴<http://cran.r-project.org/web/views/>

⁵<http://www.bioconductor.org/packages/devel/BiocViews.html>

85 2.2. Getting suitable data

86 Software development, evaluation and demonstration can not be envisioned without ap-
87 propriate data. Although R packages most often focus on software functionality, packages are
88 also used to distribute experimental and annotation data, displayed in the *AnnotationData*
89 and *ExperimentData biocViews*. A specific `MassSpectrometryData` category, currently of-
90 fering 5 packages, is dedicated for experimental data of interest here. Software packages
91 often also distribute small data sets for illustration, demonstration and code testing.

92 To exemplify some of the pipelines in this publication, we will make use of a larger,
93 public data set, available from the ProteomeXchange⁶ [22] ProteomeCentral repository (data
94 PXD000001⁷). In this TMT 6-plex [23] experiment, four exogenous proteins were spiked
95 into an equimolar *Erwinia carotovora* lysate with varying proportions in each channel of
96 quantitation; yeast enolase (ENO) at 10:5:2.5:1:2.5:10, bovine serum albumin (BSA) at
97 1:2.5:5:10:5:1, rabbit glycogen phosphorylase (PHO) at 2:2:2:2:1:1 and bovin cytochrome C
98 (CYT) at 1:1:1:1:1:2. Proteins were then digested, differentially labelled with TMT reagents,
99 fractionated by reverse phase nanoflow UPLC (nanoACQUITY, Waters), and analysed on
100 an LTQ Orbitrap Velos mass spectrometer (Thermo Scientific). Files in multiple format
101 will be used to illustrate the input/output capabilities that are available to the proteomics
102 audience. The companion package provides dedicated functions to directly download the
103 data.

104 2.3. Proteomics standards and MS data input-output

105 Proteomics is a very diverse field in terms of applications, experimental designs and file
106 formats. When dealing with a wide range of data, flexibility is often key; this is particularly
107 relevant for the R environment, which can be used for many different purposes and data
108 types. Raw mass spectrometry data comes in many different formats. While closed vendor-
109 specific binary formats are less interesting due to their limited scope, several research groups
110 as well as the HUPO Proteomics Standards Initiative (PSI) have developed open XML-based

⁶<http://www.proteomexchange.org/>

⁷Data DOI: <http://dx.doi.org/10.6019/PXD000001>

111 standards, formats and libraries to facilitate the development of vendor-agnostic tools and
112 analysis pipeline. This functionality is available through the `mzR` package [24, 25], that
113 provides a unified interface to the `mzData` [26], `mzXML` [27], `mzML` [28] as well as `netCDF`
114 formats. The `openMSfile` function opens a connection to any of these file types and enables
115 to query instrument information and raw data in a consistent way. It is generally used by
116 experienced users or developers who require maximal flexibility. For instance, `mzR` is used
117 by `xcms` [29, 30], `TargetSearch` [31] and `MSnbase` [32] for interaction with raw data.

118 Other packages provide higher level interfaces to raw data, modelled as computational
119 data containers that store data and meta-data while assuring internal coherence. Such
120 *classes* come with a set of associated *methods*, that allow the application of predefined
121 actions on class instances, also called *objects*, such as accessing specific pieces of information,
122 modifying parts of the data or producing relevant graphical representation of the data. The
123 `MSnExp` or `xcmsRaw` classes, defined in the `MSnbase` and `xcms` packages respectively, represent
124 experiments as a collection of annotated spectra, with the aim of removing the burden of
125 users to manipulate the complex data by bundling it in specialised classes with an easy-to-
126 use and well documented interface, the associated methods, to streamline the most common
127 tasks. The example raw file used below, available from the `MSnbase` package, is an iTRAQ
128 4-plex [33] experiment. It is read into R and converted into an `MSnExp` object using the
129 `readMSData` function. This specific data structure allows the spectra to be stored along
130 with associated meta data and enables easy manipulation of the complete annotated data
131 set. The last line displays a summary of the data in the R console and figure 1 illustrates
132 some of the raw data plotting functionality applicable to an `MSnExp` instance (left) or an
133 individual spectrum (right).

134 This first command finds the location of the test data file.

```
> mzXML <- dir(system.file(package = "MSnbase", dir = "extdata"),  
+ full.name = TRUE, pattern = "mzXML$")
```

135 We then proceed by reading the `mzXML` file and create an `MSnExp` object.

```
> rawms <- readMSData(mzXML, verbose = FALSE)
```

136 Finally, we show a summary of the contents of the data object.

```
> rawms

Object of class "MSnExp"
Object size in memory: 0.2 Mb
- - - Spectra data - - -
MS level(s): 2
Number of MS1 acquisitions: 1
Number of MSn scans: 5
Number of precursor ions: 5
4 unique MZs
Precursor MZ's: 437.8 - 716.34
MSn M/Z range: 100 2017
MSn retention times: 25:1 - 25:2 minutes
- - - Processing information - - -
Data loaded: Tue Apr 9 22:10:44 2013
MSnbase version: 1.9.1
- - - Meta data - - -
phenoData
  rowNames: 1
  varLabels: sampleNames fileNumbers
  varMetadata: labelDescription
Loaded from:
  dummyiTRAQ.mzXML
protocolData: none
featureData
  featureNames: X1.1 X2.1 ... X5.1 (5 total)
  fvarLabels: spectrum
  fvarMetadata: labelDescription
experimentData: use 'experimentData(object)'
```

137 [Fig. 1 about here.]

138 The `mgf` file format is also supported, for reading through the function `readMgfData`,
139 which encapsulates the peak list data into `MSnExp` objects as above, and for writing such
140 objects to a file through the `writeMgfData`. Other input/output facilities for quantified
141 data will be presented in the next section.

142 Standard formats for identification data are not yet systematically supported. It is
143 however possible to import such information into R , using existing R data import/export
144 infrastructure. For example, the XML package [34] allows one to parse arbitrary xml files
145 based on their schema definition. Support for mzIdentML, mzQuantML and possible other
146 community supported formats will be added to the mzR package.

147 *2.4. Data processing and quantitation*

148 Quantitation has become an essential part of proteomics, and several alternatives are
149 available in R for label-free and labelled approaches. In this section, we will present quanti-
150 tation functionality and associated raw data processing capabilities.

151 *2.4.1. Label-free quantitation*

152 Several packages provide functionality that can be applied to the analysis of label-free
153 MS data. Although its first scope is the study of metabolites, xcms is a mature package that
154 provides a complete pipeline for preprocessing LC/MS data for relative quantitation and
155 data visualisation [35, 36]. A typical xcms work flow implements peak extraction, filtering,
156 retention time correction and matching across samples. The package is very versatile, featur-
157 ing, for example, several peak picking methods, including some applying continuous wavelet
158 transformation (CWT) [37, 38]. The pipeline offers a complete framework to support data
159 analysis and visualisation of chromatograms and peaks to be deemed to be differentially
160 expressed. On-line help is available through a dedicated forum⁸.

161 MALDIquant [39] also provides a complete analysis pipeline for MALDI-TOF and other
162 label-free MS data. Its distinctive features include baseline subtraction using the SNIP
163 algorithm [40], peak alignment using warping functions, handling of replicated measure-
164 ments as well as supporting spectra with different resolutions. Figure 2 illustrates spectrum
165 preprocessing and peak detection steps.

166 [Fig. 2 about here.]

⁸<http://metabolomics-forum.com/>

167 `synapter` is a package [41] dedicated to the re-analysis of data independent MS^E data
168 [42, 43], acquired on Waters Synapt instruments . It implements robust data filtering strate-
169 gies, calculating and using peptide identification reliability statistics, peptide-to-protein am-
170 biguity and mass accuracy. It then models retention time deviations between reliable sets
171 of peptides in different runs and transfer identification across acquisitions to increase the
172 overall peptide and protein coverage in full experiments through an easy-to-use interface.
173 As illustrated in section 2.6, it interoperates well with `MSnbase` to take advantage of the
174 existing data structure and offers a complete analysis pipeline.

175 Finally, packages that implement MS² data processing, like `MSnbase` and `isobar` [44] (see
176 section 2.4.2), also support spectral counting once identification data is available. In addi-
177 tion, `isobar` allows one to perform emPAI [45] and distributed normalised spectral abundance
178 factor (dNSAF) [46] quantitation.

179 2.4.2. Labelled quantitation

180 Pipelines for labelled MS² quantitation, using isobaric tagging reagents such as iTRAQ
181 and TMT are available in the `isobar` and `MSnbase` packages. The code chunk below, taken
182 from `MSnbase`, illustrates how to quantify the iTRAQ reporter peaks from the `rawms` data
183 instance read in section 2.3. The `quantify` function returns another data container, an
184 `MSnSet`, specialised for storing quantitative data and associated meta data. Reporter impu-
185 rity correction can then be applied using the `purityCorrect`. The `isobar` package imports
186 centroided peak data identification data from `mgf` and text spread sheet files or converts
187 `MSnSet` instances to create its own `IBSpectra` containers for further isotope impurity cor-
188 rection, normalisation and differential expression analysis (section 2.6).

189 Below, we perform quantitation of the raw `MSnExp` data using the iTRAQ 4-plex reporters
190 ions to create a new `MSnSet` object containing the quantitative data.

```
> qnt <- quantify(rawms, reporters = iTRAQ4, verbose = FALSE)
```

191 In the following code chunk, we first define the reporter tag impurities as reporter by the
192 manufacturer, apply the correction and display a summary of the resulting `MSnSet` instance.

200 ing proteomics and metabolomics data dissemination to a wider audience through familiar
201 spreadsheet-based format, can also be incorporated and exported using the `readMzTabData`
202 and `writeMzTabData` functions. It is of course also possible to import quantitation data ex-
203 ported by third party applications to spread sheet formats. The most general way to import
204 such data is using the `read.table` function. Specialised alternatives exist, to produce data
205 structures, like `MSnSets`. The `readMSnSet` function, for instance, can import quantitation
206 data, feature meta data and sample annotation from spread sheets and create fully-fledged
207 `MSnSet` instances.

208 Additional packages provide specialised functionalities relevant to data processing. IPPD
209 [47] uses template matching to deconvolute peak patterns in individual raw spectra or com-
210 plete experiments. `Rdisop` [48, 49] is designed to determine the formula of ions based on
211 their exact mass or isotope pattern and can, reciprocally, estimate these from a formula.
212 `OrgMassSpecR` [50] has similar capabilities including specific functions to process peptide
213 and protein data: it allows the user, for example, to digest proteins, fragment peptides and
214 estimate peptide isotopic distributions modified peptides with, for example, variable ^{15}N
215 incorporation rates. In the `RforProteomics` documentation, we demonstrate how to assess
216 protein abundance of the yeast enolase spike present across the 6 PXD000001 channels using
217 `OrgMassSpecR`'s `Digest` function and observe that, allowing for one missed cleavage, we
218 observe 13 out of 79 peptides with length greater than 7 residues (corresponding to the
219 shortest identified ENO peptide), as illustrated in figure 4. The \LaTeX code producing the
220 alignment for the figure has been generated automatically, from within R, using the protein
221 sequence and observed peptide sequences and `TEXshade` [51].

222 [Fig. 4 about here.]

223 2.5. Quality control

224 Data quality is a concern in any experimental science, but the high throughput nature
225 of modern *omics* technologies, including proteomics [52, 53], requires the development of
226 specific data exploration techniques to highlight specific patterns in data. Examination of

227 complex data is greatly facilitated by well structured containers such as those cited above,
228 that enable direct access to a specific set of values. This, in turn, streamlines the implemen-
229 tation of default and robust pipelines that recurrently query the same data to produce the
230 diagnostic plots and metrics. It is however also often necessary to manually explore data
231 specificity, making the availability of data management facilities even more important.

232 In this section, we present 3 quality plots (figure 5) that can be used to assess the intrinsic
233 features of the PXD000001 data set at different levels. On the left, the distribution of MS²
234 delta m/z [54] allows the user to assess the relevance of peptide identification; high quality
235 data show m/z differences corresponding to amino acid residue masses rising well above the
236 general noise level in the histogram. One can also observe a peak at 44 Da, corresponding
237 to the mass of a polyethylene glycol (PEG) monomer, a common laboratory contaminant in
238 MS. The middle figure illustrates incomplete dissociation of TMT reporter tags, a technical
239 characteristic of the labelling approach. Incomplete dissociation of the reporter and balance
240 moieties of isobaric tags result in this additional single fragment ion peak, in which the
241 multiple channels of quantitation remain convoluted. The figure illustrates the sum of
242 genuine reporter peaks as a function of incompletely dissociated reporter data. The dotted
243 line corresponds to equal real and lost signal. A linear model has been fitted to the data
244 (blue line), indicating that there is, on average, 100-fold more genuine reporter signal. The
245 heatmap on the right indicates the relevance of our quantitation data at the level of our
246 experiment. Congruent peptide clustering indicates agreement between spike peptides while
247 no significant grouping is detected for the samples.

248 [Fig. 5 about here.]

249 Although the figures above are helpful individually, quality assessment is often most
250 efficient when put into context. Lab-wide monitoring of quality properties and metrics over
251 time to gain experience of average performances and critical thresholds, is the most efficient
252 and valuable application of quality control; the tools presented in this section are one way
253 to automate such a process.

254 *2.6. Data analysis*

255 In this section, we will describe data analysis pipelines for two quantitative strategies,
256 namely MS^Elabel-free and isobaric tagging, using **synapter** and **isobar** respectively.

257 Once quantitation data is obtained, it is often desirable to correct technical biases to
258 improve detection of biologically relevant proteins. The availability of well established nor-
259 malisation algorithms within the Bioconductor project are directly applicable here. The
260 **MSnSet** object called **qnt**, created in section 2.4.2 can be normalised using various meth-
261 ods, including quantile normalisation [55] and variance stabilisation [56, 57] using a single
262 **normalize** command. **isobar** also has similar functionality, tailored for **IBSpectra** objects;
263 its **normalize** method corrects by a factor such that the median intensities in all reporter
264 channels are equal.

265 **isobar** implements methodology to model variability in the data. We will illustrate this
266 using the PXD000001 data to estimate spectra and proteins exhibiting significant differences
267 between channel 127 and 129. As shown on figure 6, experimental noise has been approxi-
268 mated using the **NoiseModel** function on *Erwinia* background (red), spiked-in (blue) or all
269 (green) peptides (left) and protein ratios and significance have been computed (using the
270 full noise model) with the **estimateRatio** function, to call statistically relevant proteins.

271 [Fig. 6 about here.]

272 Data independent MS^E acquisition from a Synapt mass spectrometer (Waters) can be
273 efficiently analysed in R using the **synapter** pipeline, providing a complete and open work
274 flow (figure 7) leading to comprehensive data exploration and more reliable results. The test
275 data used for this illustration is a spiked-in set distributed with the **synapterdata** package: 3
276 replicates (labelled *a* to *c*) of the Universal Proteomics Standard (UPS1, Sigma) 48 protein
277 mix at 25 fmol and 3 replicates at 50 fmol, in a constant *Escherichia coli* background. The
278 set of functions in **synapter** produce data in a specific data container, called **Synapter** objects,
279 and labelled **ups** on figure 7. They store quantitative data for a set of *m* identified peptides
280 for one unique sample. Although at this step, much has been gained in terms of reliability

281 and number of peptides, we are still far from having interpretable results at this stage. These
282 **Synapter** objects can easily be converted into **MSnSet** instances (of dimensions $m_i \times 1$, where
283 m_i is the number of peptides for the processed sample, labelled **ms** on figure 7). Each newly
284 converted MS^E data can now be quantified using the top 3 method [42] (or any top n
285 variant) where the intensities of the 3 most intense peptides for each protein are aggregated
286 to estimate protein quantities. Each set of replicates is then combined into two new $m_i \times 3$
287 **MSnSet** instances (named **ms25** and **ms50**), one for each set of spike concentration, that are
288 then filtered for missing quantitation, keeping only proteins that have been quantified in at
289 least 2 out of 3 replicates. **ms25** and **ms50** are finally combined into the final $m_i \times 6$ final
290 data, normalised and subjected to a statistical analysis. As illustrated above, it becomes
291 possible to design specific pipelines for any type of experiments using standardised methods
292 and data structures.

293 [Fig. 7 about here.]

294 2.7. MS^2 spectra identification

295 A very recent addition to Bioconductor is the **rTANDEM** package [58]. The package en-
296 capsulates the mass spectrometry identification algorithm X!Tandem [59], the software for
297 protein identification by tandem mass spectrometry, in R, making it possible to perform MS^2
298 spectra identification within the R environment and directly benefit from R's data mining ca-
299 pabilities to explore the results. The package includes the X!Tandem source code eliminating
300 independent installation of the search engine. In its most basic form, the package allows to
301 call the **tandem(input)** function, where **input** is either an object of a dedicated class or the
302 path to a parameter file, as one would execute **tandem.exe /path/to/input.xml** from the
303 command line. The results are, as in the original X!Tandem software, stored in an **xml**, which
304 can however be imported into R in a straightforward way using the **GetResultsFromXML**
305 function to subsequently extract the identified peptides and inferred proteins.

306 **rTANDEM** is currently the only direct R interface to a search engine and is as such of
307 particularly noteworthy. Other alternatives require to execute the spectra identification

308 outside of R and import, export it in an appropriate format and subsequently import is into
309 R .

310 *2.8. Annotation infrastructure*

311 The Bioconductor project provides extensive annotation resources through curated off-
312 line annotation packages, that are updated with every release, or through packages that
313 provide direct on-line access to web-based repositories. The former can be targeted towards
314 specific organisms (e.g. `org.Hs.eg.db` [60] for *Homo sapiens*) of systems-level annotation
315 such as gene ontology (the `GO.db` package [61] to gain access to the Gene Ontology [62]
316 annotation) or gene pathways (the `reactome.db` [63] interface to the reactome database [64,
317 65]). `biomaRt` [66, 67] is a very flexible solution to build elaborated web queries to dedicated
318 data mart servers. Both approaches have advantages. While on-line queries allow one to
319 obtain the latest up-to-date information, they rely on network availability and immediate
320 reproducibility in less straightforward to control.

321 In the `RforProteomics` documentation, we demonstrate a use case applying 3 complemen-
322 tary alternatives. If one wishes, for example, to extract sub-cellular localisation for a gene
323 of interest, say the human HECW1 gene with Ensembl id ENSG00000002746, it is possible
324 to use (1) the `hpar` package [68] to query the Human Protein Atlas data [69, 70] or (2) to
325 query the `org.Hs.eg.db` and `GO.db` annotations to extract the relevant information or (3)
326 `biomaRt` to query the Ensembl server. Each alternative reports the same location, namely
327 nucleus and cytoplasm, although this might not be necessarily the case. The `hpar` results
328 are very specific and manually annotated, specifying that the protein, although observed in
329 the nucleus, has not been observed in the nucleoli. The other generic alternatives provide
330 additional information, including GO evidence codes.

331 To conclude this section, we also refer readers to the `rols` package [71], which provides
332 on-line access to 85 ontologies through the ontology look-up service [72, 73]. Among those
333 are the PRIDE, PSI-MS (Mass Spectrometry), PSI-MI (Molecular Interaction) PSI-MOD
334 (Protein Modifications), PSI-PAR (Protein Affinity Reagents) and PRO (Protein Ontology)
335 controlled vocabularies to name those specific to proteomics and mass spectrometry.

336 3. Conclusions

337 We have illustrated data processing and analysis on a set of test and small size data.
338 While real life data sets can be processed on commodity hardware or small servers (see
339 supplementary file of [32] and the `MSnbase-demo` vignette for reports), the sophistication of
340 the biological questions of interest and the increase in throughput of instruments requires
341 software tools to adapt and scale up. R is an interpreted language (although support for
342 byte code compilation is available through the `compiler` package) and relies in many aspects
343 on a pass-by-value semantics, slowing execution of code compared to compiled languages
344 and pass-by-reference semantics. Fortunately, R's ability to interoperate with many other
345 languages, including C and C++ [74], allows users to execute computationally demanding
346 tasks while still retaining the flexibility and interactivity of the R environment. Direct
347 support for parallel computing, large memory/out-of-memory data (see for instance High-
348 Performance Computing task view¹⁰) and cloud deployment with the Bioconductor Amazon
349 Machine Image¹¹, make it possible to embark on large-scale data processing tasks.

350 Among the brief list of packages that has been reviewed, we have demonstrated alter-
351 native and complementary functionality. Most noteworthy however, is the interoperability
352 of these packages, as illustrated in some of the examples. Generally, no specific effort is
353 expected from developers to explicitly promote interaction among packages (on CRAN for
354 example), and thus it is often the user's/programmer's responsibility to implement interoper-
355 ability. The Bioconductor project, on the other hand, openly promotes interoperability
356 between packages and reuse of existing infrastructure. The classes for raw and processed
357 data, briefly described in sections 2.3 and 2.4 are adapted from and compatible with ex-
358 isting implementations for transcriptomics data, widely used in many core Bioconductor
359 packages. Data processing procedures used for data normalisation and statistical algorithms
360 are a direct and invaluable side effects of the R language and previous Bioconductor devel-
361 opment. The quality and diversity of available software, fostered by interdisciplinary, open

¹⁰<http://cran.r-project.org/web/views/HighPerformanceComputing.html>

¹¹<http://bioconductor.org/help/bioconductor-cloud-ami/>

362 and distributed development, is an immense source of knowledge to build upon.

363 Although an elaborated environment and programming language like R has undeniable
364 strengths, its sheer power and flexibility is its Achilles' heel. An important obstacle in the
365 adoption of R is its command line interface (CLI) that a user needs to apprehend before being
366 able to fully appreciate R. Life scientists very often expect to operate a software through a
367 graphical user interface (GUI), which is probably the major hurdle to the wider adoption
368 of R, or other command line environments, outside the bioinformatics community. The
369 important point is, however, that properly designed graphical and command-line interfaces
370 are good at different tasks. Flexibility, programmability and reproducibility are the strength
371 of the latter, while interactivity and navigability are the main features of the former and
372 these respective advantages are complementary. Users should not be misguided and adhere
373 to any interface through dogma or ignorance, but choose the best suited tools for any task
374 to tackle the real difficulty, which is the underlying biology.

375 In this review, we have described how to use R and a selection of packages to analyse
376 mass spectrometry based proteomics data, ranging from raw data access and visualisation,
377 data processing, labelled and label-free quantitation, quality control and data analysis. It is
378 however essential to underline that, beyond the utilisation of the functionality exposed by
379 the software, fundamental principles of data analysis have been demonstrated.

380 Every use case that is summarised, including generation of the figures, is documented
381 in the `RforProteomics` package and is fully *reproducible*: we provide code and data so that
382 interested readers are in a position to repeat the exact same steps and reproduce the same
383 results. The complexity of biological data itself and the processing it undergoes make it
384 very difficult, even for experienced users, to track the computations and verify the results by
385 merely looking at the input and the output data. As such, *transparency* of the pipeline is a
386 required condition to aim for robustness and validity of the work flow, and the software itself.
387 Biology is, by nature, extremely diverse, and creativity in the designs of experiments and the
388 development and application of technology is the main obstacle to our understanding. The
389 software that is employed must be *flexible* and extensible, to support researchers in their

390 quest rather than limit and constrain them. Reproducibility, transparency and flexibility
391 are essential characteristics for scientific software, that are provided by the tools described
392 above.

393 Despite these indisputable advantages, a lot of work still needs to be done to improve
394 and integrate our pipelines, demonstrate how R can efficiently, reproducibly and robustly be
395 used for in-depth proteomics data comprehension as well as broaden access to these tools
396 to the proteomics community. The RforProteomics is one effort in that direction. Finally,
397 support is an essential part of the success and adoption of software; the on-line R community
398 in general and the the Bioconductor mailing lists¹² in particular are a rich and broad source
399 of information for new and experienced users.

400 Acknowledgement

401 The authors are grateful to the R and Bioconductor communities for providing quality
402 software, robust data analysis methodology and helpful support. This work was supported
403 by the PRIME-XS project, grant agreement number 262067, funded by the European Union
404 7th Framework Program.

405 References

- 406 [1] T. Nilsson, M. Mann, R. Aebersold, J. R. Yates, A. Bairoch, J. J. M. Bergeron, Mass spectrometry in
407 high-throughput proteomics: ready for the big time., *Nat. Methods* 7 (9) (2010) 681–5.
- 408 [2] R. Aebersold, Editorial: From data to results, *Molecular & Cellular Proteomics* 10 (11).
- 409 [3] F. F. Gonzalez-Galarza, C. Lawless, S. J. Hubbard, J. Fan, C. Bessant, H. Hermjakob, A. R. Jones, A
410 critical appraisal of techniques, software packages, and standards for quantitative proteomic analysis.,
411 *OMICS* 16 (9) (2012) 431–42.
- 412 [4] Y. Perez-Riverol, R. Wang, H. Hermjakob, V. Vesada, J. A. Vizcaíno, Software libraries for mass
413 spectrometry based proteomics: A developers perspective, *BBA – Proteins and Proteomics*.

¹²<http://bioconductor.org/help/ mailing-list/>

- 414 [5] R Core Team, [R: A Language and Environment for Statistical Computing](#), R Foundation for Statistical
415 Computing, Vienna, Austria, ISBN 3-900051-07-0 (2012).
416 URL <http://www.R-project.org/>
- 417 [6] R. Ihaka, R. Gentleman, R: A language for data analysis and graphics, *Journal of Computational and*
418 *Graphical Statistics* 5 (3) (1996) 299–314.
- 419 [7] A. Vance, [Data analysts captivated by Rs power](#), *The New York Times*.
420 URL <http://www.nytimes.com/2009/01/07/technology/business-computing/07program.html>
- 421 [8] R. Gentleman, *R Programming for Bioinformatics*, Chapman & Hall/CRC, 2008, ISBN 978-1-420-
422 06367-7.
- 423 [9] R. C. Gentleman, V. J. Carey, D. M. Bates, B. Bolstad, M. Dettling, S. Dudoit, B. Ellis, L. Gautier,
424 Y. Ge, J. Gentry, K. Hornik, T. Hothorn, W. Huber, S. Iacus, R. Irizarry, F. Leisch, C. Li, M. Maechler,
425 A. J. Rossini, G. Sawitzki, C. Smith, G. Smyth, L. Tierney, J. Y. H. Yang, J. Zhang, *Bioconductor:*
426 *open software development for computational biology and bioinformatics.*, *Genome Biol* 5 (10) (2004)
427 –80.
- 428 [10] R. Gentleman, D. T. Lang, *Statistical analyses and reproducible research*, *Bioconductor Project Work-*
429 *ing Papers*. Working Paper 2.
- 430 [11] R. Gentleman, *Reproducible research: A bioinformatics case study*, *Statistical Applications in Genetics*
431 *and Molecular Biology* 4 (1).
- 432 [12] R. D. Peng, *Reproducible research and biostatistics.*, *Biostatistics* 10 (3) (2009) 405–408.
- 433 [13] D. L. Donoho, *An invitation to reproducible computational research.*, *Biostatistics* 11 (3) (2010) 385–8.
- 434 [14] R. D. Peng, *Reproducible research in computational science.*, *Science* 334 (6060) (2011) 1226–1227.
- 435 [15] D. E. Knuth, *Literate programming*, *The Computer Journal (British Computer Society)* 27 (2) (1984)
436 91–111.
- 437 [16] F. Leisch, *Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis*, in:
438 W. Härdle, B. Rönz (Eds.), *Compstat 2002, Proceedings in Computational Statistics*, Physica Ver-
439 lag, Heidelberg, Germany, 2002.
- 440 [17] Y. Xie, [knitr: A general-purpose package for dynamic report generation in R](#), *r* package version 1.0.5
441 (2013).
442 URL <http://CRAN.R-project.org/package=knitr>
- 443 [18] J. M. Chambers, *Software for Data Analysis: Programming with R*, Springer, New York, 2008.
- 444 [19] D. G. Messerschmitt, C. Szyperski, *Software Ecosystem: Understanding an Indispensable Technology*
445 *and Industry*, MIT Press, Cambridge, MA, USA, 2003.
- 446 [20] M. Lungu, *Reverse engineering software ecosystems*, Ph.D. thesis, University of Lugano (2009).
- 447 [21] J. Fox, *Aspects of the Social Organization and Trajectory of the R Project*, *The R Journal* 1 (2) (2009)

448 5–13.

- 449 [22] H. Hermjakob, R. Apweiler, The proteomics identifications database (pride) and the proteomexchange
450 consortium: making proteomics data accessible., *Expert Rev Proteomics* 3 (1) (2006) 1–3.
- 451 [23] A. Thompson, J. Schäfer, K. Kuhn, S. Kienle, J. Schwarz, G. Schmidt, T. Neumann, R. Johnstone,
452 A. K. A. Mohammed, C. Hamon, Tandem mass tags: a novel quantification strategy for comparative
453 analysis of complex protein mixtures by MS/MS., *Anal. Chem.* 75 (8) (2003) 1895–904.
- 454 [24] B. Fischer, S. Neumann, L. Gatto, [mzR: parser for netCDF, mzXML, mzData and mzML files \(mass
455 spectrometry data\)](#), R package version 1.3.9 (2012).
456 URL <http://www.bioconductor.org/packages/release/bioc/html/mzR.html>
- 457 [25] M. C. Chambers, B. Maclean, R. Burke, D. Amodei, D. L. Ruderman, S. Neumann, L. Gatto, B. Fischer,
458 B. Pratt, J. Egertson, K. Hoff, D. Kessner, N. Tasman, N. Shulman, B. Frewen, T. A. Baker, M. Y.
459 Brusniak, C. Paulse, D. Creasy, L. Flashner, K. Kani, C. Moulding, S. L. Seymour, L. M. Nuwaysir,
460 B. Lefebvre, F. Kuhlmann, J. Roark, P. Rainer, S. Detlev, T. Hemenway, A. Huhmer, J. Langridge,
461 B. Connolly, T. Chadick, K. Holly, J. Eckels, E. W. Deutsch, R. L. Moritz, J. E. Katz, D. B. Agus,
462 M. MacCoss, D. L. Tabb, P. Mallick, A cross-platform toolkit for mass spectrometry and proteomics.,
463 *Nat Biotechnol* 30 (10) (2012) 918–20. doi:10.1038/nbt.2377.
- 464 [26] S. Orchard, L. Montechi-Palazzi, E. W. Deutsch, P.-A. Binz, A. R. Jones, N. Paton, A. Pizarro, D. M.
465 Creasy, J. Wojcik, H. Hermjakob, Five years of progress in the standardization of proteomics data 4th
466 annual spring workshop of the hupo-proteomics standards initiative april 23-25, 2007 ecole nationale
467 supérieure (ens), lyon, france., *Proteomics* 7 (19) (2007) 3436–40.
- 468 [27] P. G. A. Pedrioli, J. K. Eng, R. Hubley, M. Vogelzang, E. W. Deutsch, B. Raught, B. Pratt, E. Nilsson,
469 R. H. Angeletti, R. Apweiler, K. Cheung, C. E. Costello, H. Hermjakob, S. Huang, R. K. Julian,
470 E. Kapp, M. E. McComb, S. G. Oliver, G. Omenn, N. W. Paton, R. Simpson, R. Smith, C. F. Taylor,
471 W. Zhu, R. Aebersold, A common open representation of mass spectrometry data and its application
472 to proteomics research., *Nat. Biotechnol.* 22 (11) (2004) 1459–66.
- 473 [28] L. Martens, M. Chambers, M. Sturm, D. Kessner, F. Levander, J. Shofstahl, W. H. Tang, A. Rompp,
474 S. Neumann, A. D. Pizarro, L. Montechi-Palazzi, N. Tasman, M. Coleman, F. Reisinger, P. Souda,
475 H. Hermjakob, P.-A. Binz, E. W. Deutsch, mzML - a community standard for mass spectrometry data.,
476 *Molecular and Cellular Proteomics*.
- 477 [29] C. A. Smith, E. J. Want, G. O’Maille, R. Abagyan, G. Siuzdak, XCMS: processing mass spectrometry
478 data for metabolite profiling using nonlinear peak alignment, matching, and identification., *Anal Chem*
479 78 (3) (2006) 779–87.
- 480 [30] H. P. Benton, D. M. Wong, S. A. Trauger, G. Siuzdak, XCMS2: processing tandem mass spectrometry
481 data for metabolite identification and structural characterization., *Anal Chem* 80 (16) (2008) 6382–9.

- 482 [31] A. Cuadros-Inostroza, C. Caldana, H. Redestig, J. Lisec, H. Pena-Cortes, L. Willmitzer, M. A. Hannah,
483 TargetSearch - a Bioconductor package for the efficient pre-processing of GC-MS metabolite profiling
484 data, *BMC Bioinformatics* 10 (2009) 428.
- 485 [32] L. Gatto, K. S. Lilley, MSnbase – an R/Bioconductor package for isobaric tagged mass spectrometry
486 data visualization, processing and quantitation., *Bioinformatics* 28 (2) (2012) 288–9.
- 487 [33] P. L. Ross, Y. N. Huang, J. N. Marchese, B. Williamson, K. Parker, S. Hattan, N. Khainovski, S. Pillai,
488 S. Dey, S. Daniels, S. Purkayastha, P. Juhasz, S. Martin, M. Bartlet-Jones, F. He, A. Jacobson, D. J.
489 Pappin, Multiplexed protein quantitation in *Saccharomyces cerevisiae* using amine-reactive isobaric
490 tagging reagents., *Mol Cell Proteomics* 3 (12) (2004) 1154–1169.
- 491 [34] D. T. Lang, [XML: Tools for parsing and generating XML within R and S-Plus.](#), R package version
492 3.9-4 (2012).
493 URL <http://CRAN.R-project.org/package=XML>
- 494 [35] L. N. Mueller, M. Y. Brusniak, D. R. Mani, R. Aebersold, An assessment of software solutions for the
495 analysis of mass spectrometry based quantitative proteomics data., *J Proteome Res* 7 (1) (2008) 51–61.
- 496 [36] E. Lange, R. Tautenhahn, S. Neumann, C. Grpl, Critical assessment of alignment procedures for LC-MS
497 proteomics and metabolomics measurements., *BMC Bioinformatics* 9 (2008) 375.
- 498 [37] P. Du, W. A. Kibbe, S. M. Lin, Improved peak detection in mass spectrum by incorporating continuous
499 wavelet transform-based pattern matching., *Bioinformatics* 22 (17) (2006) 2059–65.
- 500 [38] R. Tautenhahn, C. Böttcher, S. Neumann, Highly sensitive feature detection for high resolution
501 LC/MS., *BMC Bioinformatics* 9 (2008) 504.
- 502 [39] S. Gibb, K. Strimmer, MALDIquant: a versatile R package for the analysis of mass spectrometry data.,
503 *Bioinformatics* 28 (17) (2012) 2270–1.
- 504 [40] C. G. Ryan, E. Clayton, W. L. Griffin, S. H. Sie, D. R. Cousens, SNIP, a statistics-sensitive background
505 treatment for the quantitative analysis of PIXE spectra in geoscience applications, *Nuclear Instruments*
506 *and Methods in Physics Research B* 34 (1988) 396–402.
- 507 [41] L. Gatto, P. V. Shliaha, N. J. Bond, [synapter: Label-free data analysis pipeline for optimal identification](#)
508 [and quantitation](#), R package version 0.99.13 (2012).
509 URL <http://bioconductor.org/packages/devel/bioc/html/synapter.html>
- 510 [42] J. C. Silva, M. V. Gorenstein, G. Z. Li, J. P. Vissers, S. J. Geromanos, Absolute quantification of
511 proteins by lcms: a virtue of parallel ms acquisition., *Mol Cell Proteomics* 5 (1) (2006) 144–56.
- 512 [43] S. J. Geromanos, J. P. Vissers, J. C. Silva, C. A. Dorschel, G. Z. Li, M. V. Gorenstein, R. H. Bateman,
513 J. I. Langridge, The detection, correlation, and comparison of peptide precursor and product ions from
514 data independent LC-MS with data dependant LC-MS/MS., *Proteomics* 9 (6) (2009) 1683–95.
- 515 [44] F. P. Breitwieser, A. Mller, L. Dayon, T. Kcher, A. Hainard, P. Pichler, U. Schmidt-Erfurth, G. Superti-

- 516 Furga, J. C. Sanchez, K. Mechtler, K. L. Bennett, J. Colinge, General statistical modeling of data from
517 protein relative expression isobaric tags., *J Proteome Res* 10 (6) (2011) 2758–66.
- 518 [45] Y. Ishihama, Y. Oda, T. Tabata, T. Sato, T. Nagasu, J. Rappsilber, M. Mann, Exponentially modified
519 protein abundance index (emPAI) for estimation of absolute protein amount in proteomics by the
520 number of sequenced peptides per protein., *Mol Cell Proteomics* 4 (9) (2005) 1265–1272.
- 521 [46] Y. Zhang, Z. Wen, M. P. Washburn, L. Florens, Refinements to label free proteome quantitation: how
522 to deal with peptides shared by multiple proteins., *Anal Chem* 82 (6) (2010) 2272–81.
- 523 [47] M. Slawski, R. Hussong, M. Hein, *IPPD: Isotopic peak pattern deconvolution for Protein Mass Spec-*
524 *trometry by template matching*, R package version 1.5.0 (2012).
525 URL <http://www.bioconductor.org/packages/release/bioc/html/IPPD.html>
- 526 [48] S. Böcker, Z. Lipták, M. Martin, A. Pervukhin, H. Sudek, DECOMP – from interpreting mass spec-
527 trometry peaks to solving the money changing problem., *Bioinformatics* 24 (4) (2008) 591–3.
- 528 [49] S. Böcker, M. C. Letzel, Z. Lipták, A. Pervukhin, SIRIUS: decomposing isotope patterns for metabolite
529 identification., *Bioinformatics* 25 (2) (2009) 218–24.
- 530 [50] N. G. Dodder, with code contributions from Katharine M. Mullen., *OrgMassSpecR: Organic Mass*
531 *Spectrometry*, R package version 0.3-12 (2012).
532 URL <http://CRAN.R-project.org/package=OrgMassSpecR>
- 533 [51] E. Beitz, Texshade: shading and labeling of multiple sequence alignments using latex2e, *Bioinformatics*
534 (2000) 135–139.
- 535 [52] A. Beasley-Green, D. Bunk, P. Rudnick, L. Kilpatrick, K. Phinney, A proteomics performance standard
536 to support measurement quality in proteomics., *Proteomics* 12 (7) (2012) 923–31.
- 537 [53] Z. Q. Ma, K. O. Polzin, S. Dasari, M. C. Chambers, B. Schilling, B. W. Gibson, B. Q. Tran, L. Vega-
538 Montoto, D. C. Liebler, D. L. Tabb, QuaMeter: Multivendor performance metrics for LC-MS/MS
539 proteomics instrumentation., *Anal Chem* 84 (14) (2012) 5845–50.
- 540 [54] J. M. Foster, S. Degroeve, L. Gatto, M. Visser, R. Wang, J. Griss, R. Apweiler, L. Martens, A posteriori
541 quality control for the curation and reuse of public proteomics data., *Proteomics* 11 (11) (2011) 2182–94.
- 542 [55] B. M. Bolstad, R. A. Irizarry, M. Astrand, T. P. Speed, A comparison of normalization methods for high
543 density oligonucleotide array data based on variance and bias., *Bioinformatics* 19 (2) (2003) 185–93.
- 544 [56] W. Huber, A. von Heydebreck, H. Suelmann, A. Poustka, M. Vingron, Variance stabilization applied
545 to microarray data calibration and to the quantification of differential expression, *Bioinformatics* 18
546 Suppl. 1 (2002) S96–S104.
- 547 [57] N. A. Karp, W. Huber, P. G. Sadowski, P. D. Charles, S. V. Hester, K. S. Lilley, Addressing accuracy
548 and precision issues in iTRAQ quantitation., *Mol. Cell Proteomics* 9 (9) (2010) 1885–97.
- 549 [58] F. Fournier, C. J. Beauparlant, R. Paradis, A. Droit, rTANDEM: Encapsulate X!Tandem in R., r

550 package version 0.99.4 (2013).

551 [59] R. Craig, R. C. Beavis, Tandem: matching proteins with tandem mass spectra., *Bioinformatics* 20 (9)
552 (2004) 1466–7. doi:10.1093/bioinformatics/bth092.

553 [60] M. Carlson, [org.Hs.eg.db: Genome wide annotation for Human](http://org.Hs.eg.db), R package version 2.8.0 (2012).
554 URL <http://bioconductor.org/packages/release/data/annotation/html/org.Hs.eg.db.html>

555 [61] M. Carlson, [GO.db: A set of annotation maps describing the entire Gene Ontology](http://GO.db), R package version
556 2.8.0 (2012).
557 URL <http://bioconductor.org/packages/release/data/annotation/html/GO.db.html>

558 [62] M. Ashburner, C. A. Ball, J. A. Blake, D. Botstein, H. Butler, J. M. Cherry, A. P. Davis, K. Dolinski,
559 S. S. Dwight, J. T. Eppig, M. A. Harris, D. P. Hill, L. Issel-Tarver, A. Kasarskis, S. Lewis, J. C. Matese,
560 J. E. Richardson, M. Ringwald, G. M. Rubin, G. Sherlock, Gene ontology: tool for the unification of
561 biology. the gene ontology consortium., *Nat Genet* 25 (1) (2000) 25–9.

562 [63] W. Ligtenberg, [reactome.db: A set of annotation maps for reactome](http://reactome.db), R package version 1.1.0 (2012).
563 URL <http://bioconductor.org/packages/release/data/annotation/html/reactome.db.html>

564 [64] D. Croft, G. O’Kelly, G. Wu, R. Haw, M. Gillespie, L. Matthews, M. Caudy, P. Garapati, G. Gopinath,
565 B. Jassal, S. Jupe, I. Kalatskaya, S. Mahajan, B. May, N. Ndegwa, E. Schmidt, V. Shamovsky, C. Yung,
566 E. Birney, H. Hermjakob, P. D’Eustachio, L. Stein, Reactome: a database of reactions, pathways and
567 biological processes., *Nucleic Acids Res* 39 (Database issue) (2011) D691–7.

568 [65] P. D’Eustachio, Reactome knowledgebase of human biological pathways and processes., *Methods Mol*
569 *Biol* 694 (2011) 49–61.

570 [66] S. Durinck, Y. Moreau, A. Kasprzyk, S. Davis, B. De Moor, A. Brazma, W. Huber, Biomart and bio-
571 conductor: a powerful link between biological databases and microarray data analysis., *Bioinformatics*
572 21 (16) (2005) 3439–40.

573 [67] S. Durinck, P. T. Spellman, E. Birney, W. Huber, Mapping identifiers for the integration of genomic
574 datasets with the r/bioconductor package biomart., *Nat Protoc* 4 (8) (2009) 1184–91.

575 [68] L. Gatto, [hpar: Human Protein Atlas in R](http://hpar), R package version 0.99.0 (2012).
576 URL <http://bioconductor.org/packages/devel/bioc/html/hpar.html>

577 [69] M. Uhlén, E. Björling, C. Agaton, C. A.-K. A. Szgyarto, B. Amini, E. Andersen, A.-C. C. Andersson,
578 P. Angelidou, A. Asplund, C. Asplund, L. Berglund, K. Bergström, H. Brumer, D. Cerjan, M. Ekström,
579 A. Elobeid, C. Eriksson, L. Fagerberg, R. Falk, J. Fall, M. Forsberg, M. G. G. Björklund, K. Gumbel,
580 A. Halimi, I. Hallin, C. Hamsten, M. Hansson, M. Hedhammar, G. Hercules, C. Kampf, K. Larsson,
581 M. Lindskog, W. Lodewyckx, J. Lund, J. Lundeberg, K. Magnusson, E. Malm, P. Nilsson, J. Odling,
582 P. Oksvold, I. Olsson, E. Oster, J. Ottosson, L. Paavilainen, A. Persson, R. Rimini, J. Rockberg,
583 M. Runeson, A. Sivertsson, A. Skölleremo, J. Steen, M. Stenvall, F. Sterky, S. Strömberg, M. Sundberg,

- 584 H. Tegel, S. Tourle, E. Wahlund, A. Waldén, J. Wan, H. Wernérus, J. Westberg, K. Wester, U. Wretha-
585 gen, L. L. L. Xu, S. Hober, F. Pontén, A human protein atlas for normal and cancer tissues based on
586 antibody proteomics., *Molecular & cellular proteomics : MCP* 4 (12) (2005) 1920–1932.
- 587 [70] M. Uhlen, P. Oksvold, L. Fagerberg, E. Lundberg, K. Jonasson, M. Forsberg, M. Zwahlen, C. Kampf,
588 K. Wester, S. Hober, H. Wernerus, L. Björling, F. Ponten, Towards a knowledge-based Human Protein
589 Atlas., *Nature biotechnology* 28 (12) (2010) 1248–1250.
- 590 [71] L. Gatto, [rols: An R interface to the Ontology Lookup Service](https://rols.bioconductor.org/), R package version 0.99.10 (2012).
591 URL <http://bioconductor.org/packages/devel/bioc/html/rols.html>
- 592 [72] R. G. Côté, P. Jones, R. Apweiler, H. Hermjakob, The ontology lookup service, a lightweight cross-
593 platform tool for controlled vocabulary queries., *BMC Bioinformatics* 7 (2006) 97.
- 594 [73] R. G. Côté, P. Jones, L. Martens, R. Apweiler, H. Hermjakob, The ontology lookup service: more data
595 and better tools for controlled vocabulary queries., *Nucleic Acids Res.* 36 (Web Server issue) (2008)
596 372–376.
- 597 [74] D. Eddelbuettel, R. François, Rcpp: Seamless R and C++ integration, *Journal of Statistical Software*
598 40 (8) (2011) 1–18.
- 599 [75] H. Wickham, *ggplot2: elegant graphics for data analysis*, Springer New York, 2009.

600 **List of Figures**

601 1 Plotting raw MS² data using functionality from the **MSnbase** package. On the
602 left, the full *m/z* range of an experiment containing 5 spectra is displayed. On
603 the right, one spectrum of interest is illustrated, highlighting the 4 iTRAQ re-
604 porter region. Both figures, have been created with the generic **plot** function,
605 applied to either the complete experiment of a single MS² spectrum. 26

606 2 Label-free spectrum processing peak detection from the **MALDIquant** package.
607 Figures represent (1) raw data, (2) effect of variance stabilisation using square
608 root transformation, (3) smoothing using a simple 5 point moving average,
609 (4) base line correction, (5) noise reduction and peak detection and (6) final
610 results. 27

611 3 Representation of peptide-level quantitation data. This plot has been gener-
612 ated using the PXD000001 TMT 6-plex data and converted to an **MSnSet**
613 object. Normalised background and spike (BSA, CYT, ENO and PHO) re-
614 porter ion intensities for a subset of peptides have been plotted using the
615 **ggplot2** package [75]. The complete code is available in the companion package. 28

616 4 Visualising observed peptides for the yeast enolase protein. Consecutive pep-
617 tides are shaded in different colours. The last peptide is a miscleavage and
618 overlaps with IEEELGDNAVFAGENFHHGDK. 29

619 5 Assessing the quality of the PXD000001 data set. On the left, the delta *m/z*
620 plot illustrates the relevance of the raw MS² spectra for peptide identification.
621 The middle figure compares fully dissociated reporter signal against incom-
622 pletely dissociated ions, indicating satisfactory reporter dissociation for the
623 experiment. The last figure, a heatmap of a subset of peptides, highlights the
624 expected lack of sample grouping and tight peptides clustering. The first plot
625 is produced by the **plotMzDelta** function from the **MSnbase** package. The
626 other figures used standard base R plotting functionality. The detailed code
627 and data to reproduce the figures is available in companion package. 30

628 6 On the left, the MA plot for the PXD000001 127 vs. 129 reporter ions, showing
629 the 95% confidence intervals of the background peptides (red), spikes (blue)
630 and all (green) peptide noise models. The respective peptides are colour-coded
631 according to the proteins. The volcano plot on the right illustrates protein
632 significance ($-\log_{10}$ p-value) as a function of the \log_{10} fold-change. The ver-
633 tical coloured dashed indicate the expected \log_{10} ratios. The black dotted
634 horizontal and vertical lines represent a p-value of 0.01 and fold-changes of
635 0.5 and 2 respectively. 31

636 7 The **synapter** to **MSnbase** pipeline, illustrating how to combine and process
637 data objects in an design specific work flow. Data objects are represented by
638 grey boxes, while functions, that manipulate and transform the objects are
639 shown in white boxes. The respective dimensions of the objects (number of
640 features \times number of sample) are given in parenthesis. 32

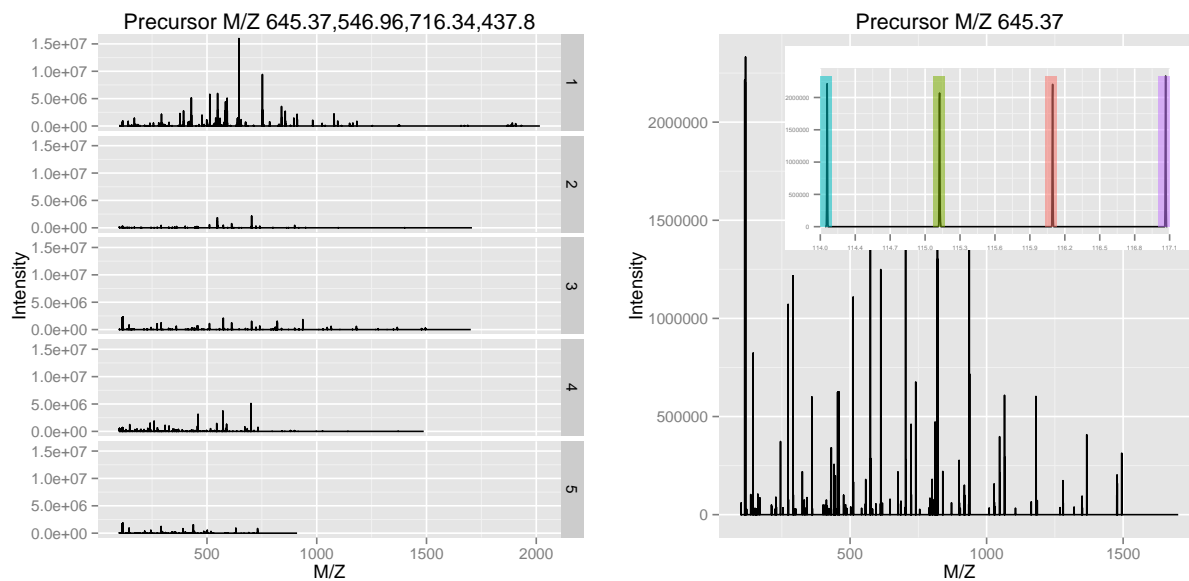


Fig. 1: Plotting raw MS² data using functionality from the MSnbase package. On the left, the full m/z range of an experiment containing 5 spectra is displayed. On the right, one spectrum of interest is illustrated, highlighting the 4 iTRAQ reporter region. Both figures, have been created with the generic `plot` function, applied to either the complete experiment or a single MS² spectrum.

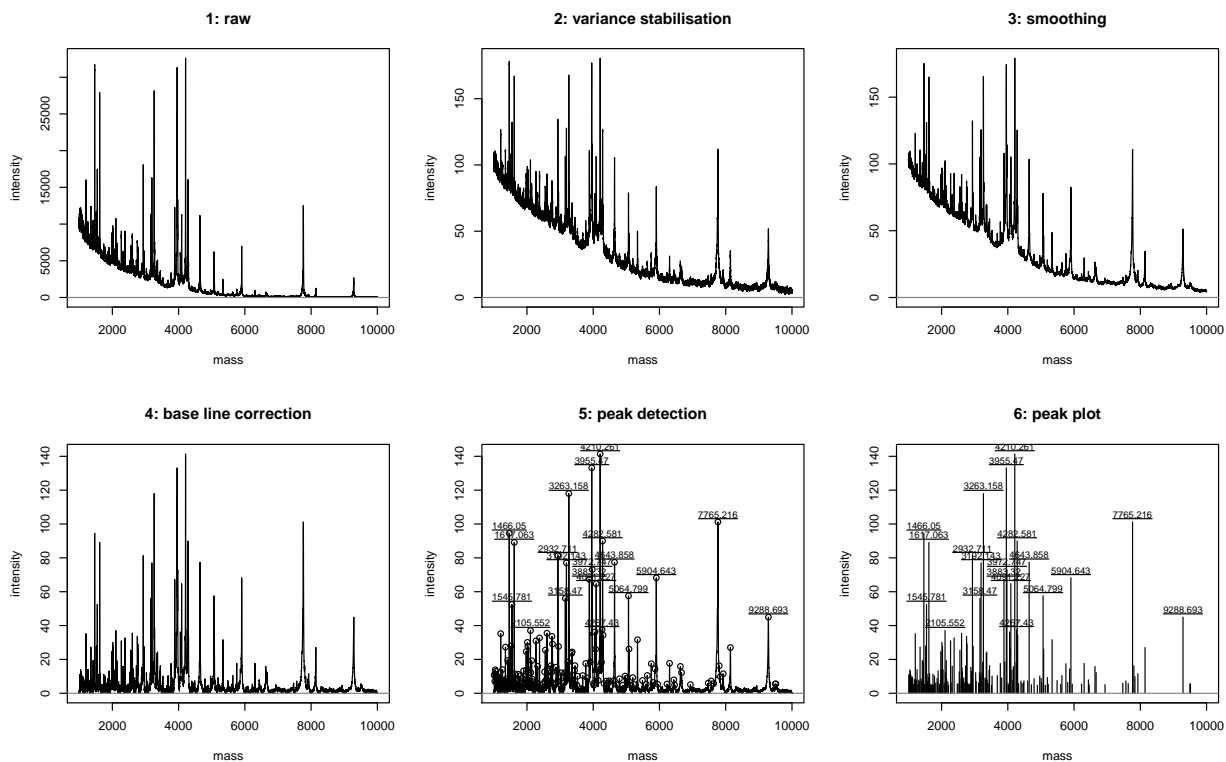


Fig. 2: Label-free spectrum processing peak detection from the MALDIquant package. Figures represent (1) raw data, (2) effect of variance stabilisation using square root transformation, (3) smoothing using a simple 5 point moving average, (4) base line correction, (5) noise reduction and peak detection and (6) final results.

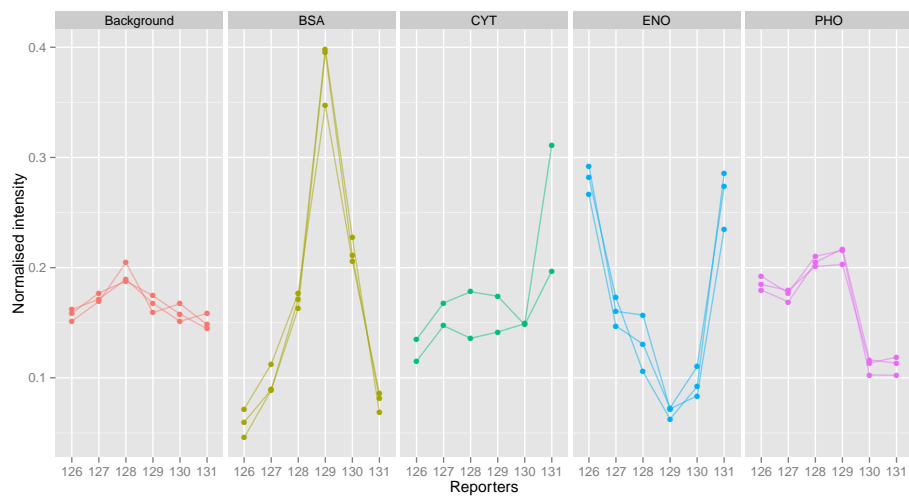


Fig. 3: Representation of peptide-level quantitation data. This plot has been generated using the PXD000001 TMT 6-plex data and converted to an MSnSet object. Normalised background and spike (BSA, CYT, ENO and PHO) reporter ion intensities for a subset of peptides have been plotted using the ggplot2 package [75]. The complete code is available in the companion package.

MAVSKVYARSVYDSRGNPTVEVELTTEKGVFRSIVPSGASTGVHEALEMRDGDKSKWMGKGVLHAVKNVN 70
 DVIAPAFVKANIDVKDQKAVDDFLISLDGTANKSKLGANAILGVSLAASRAAAAEKNVPLYKHLADLSKS 140
 KTSPYVLPVPFLNVLNGGSHAGGALALQEFMIAPTGAKTFAEALRIGSEVYHNLKSLTKKRYGASAGNVG 210
 DEGGVAPNIQTAEALDLIVDAIKAAGHDGKIKIGLDCASSEFFKDGKYDLDFKNPNSDKSKWLTGPQLA 280
 DLYHSLMKRYPIVSIEDPFAEDDWEAWSHFFKTAGIQIVADDLTVTNPKRIATAIEKKAADALLKVNQI 350
 GTLSESIKAAQDSFAAGWGMVSHRSGETEDTFIADLVVGLRTGQIKTGAPARSERLAKLNQLLRIEEEL 420
 GDNAVFAGENFHHGDKL 437

Fig. 4: Visualising observed peptides for the yeast enolase protein. Consecutive peptides are shaded in different colours. The last peptide is a miscleavage and overlaps with IEEELGDNAVFAGENFHHGDK.

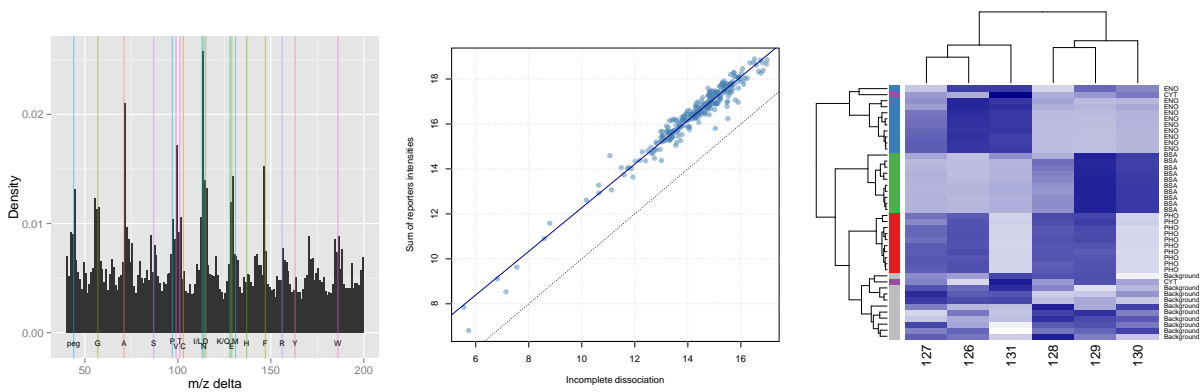


Fig. 5: Assessing the quality of the PXD000001 data set. On the left, the delta m/z plot illustrates the relevance of the raw MS² spectra for peptide identification. The middle figure compares fully dissociated reporter signal against incompletely dissociated ions, indicating satisfactory reporter dissociation for the experiment. The last figure, a heatmap of a subset of peptides, highlights the expected lack of sample grouping and tight peptides clustering. The first plot is produced by the `plotMzDelta` function from the `MSnbase` package. The other figures used standard base R plotting functionality. The detailed code and data to reproduce the figures is available in companion package.

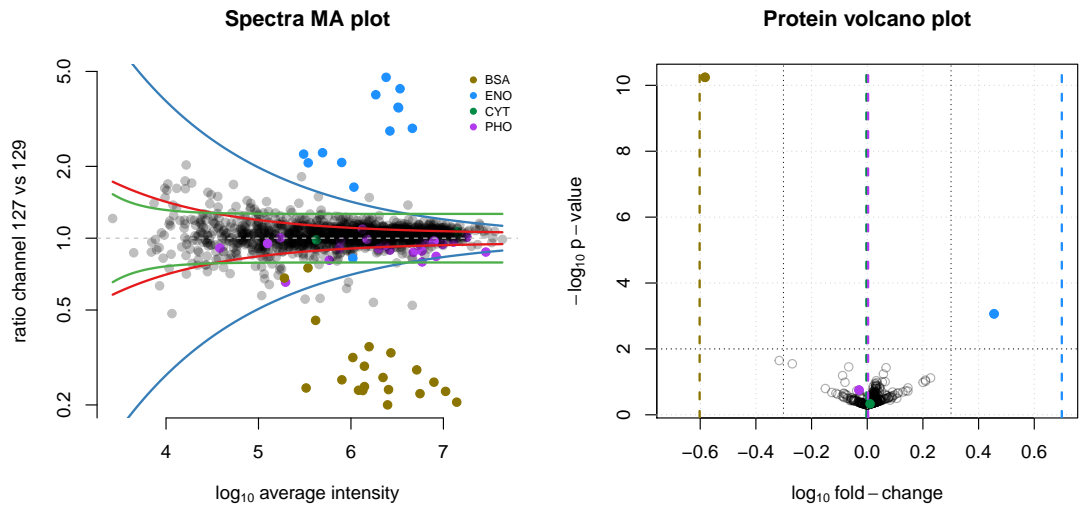


Fig. 6: On the left, the MA plot for the PXD000001 127 vs. 129 reporter ions, showing the 95% confidence intervals of the background peptides (red), spikes (blue) and all (green) peptide noise models. The respective peptides are colour-coded according to the proteins. The volcano plot on the right illustrates protein significance ($-\log_{10}$ p-value) as a function of the \log_{10} fold-change. The vertical coloured dashed indicate the expected \log_{10} ratios. The black dotted horizontal and vertical lines represent a p-value of 0.01 and fold-changes of 0.5 and 2 respectively.

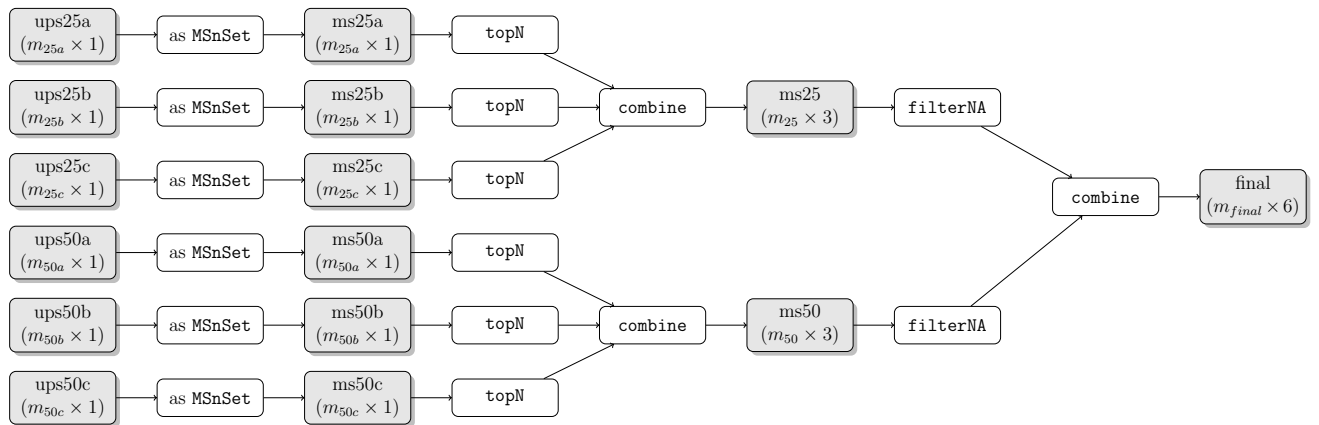


Fig. 7: The synapter to MSnbase pipeline, illustrating how to combine and process data objects in an design specific work flow. Data objects are represented by grey boxes, while functions, that manipulate and transform the objects are shown in white boxes. The respective dimensions of the objects (number of features \times number of sample) are given in parenthesis.