

1 Article / Concept Paper

2

### 3 A method to encourage minimum reporting guideline uptake for data analysis in 4 metabolomics

5 Elizabeth C. Considine <sup>1\*</sup>, Reza M. Salek <sup>2</sup>

6 <sup>1</sup>The Irish Centre for Fetal and Neonatal Translational Research (INFANT), Department of  
7 Obstetrics and Gynaecology, University College Cork, Cork, Ireland

8 <sup>2</sup> The International Agency for Research on Cancer (IARC), 150 Cours Albert Thomas,  
9 69372 Lyon CEDEX 08, France.

10 \* lizconsidine@gmail.com

11

12

#### 13 **Abstract: Introduction**

14 Despite the proposal of minimum reporting guidelines for metabolomics over a decade ago,  
15 reporting on the data analysis step in metabolomics has been shown to be unclear and  
16 incomplete with major omissions and lack of logical flow rendering the data analysis'  
17 workflows in these studies impossible to follow and therefore replicate or even imitate. Here  
18 we propose possible reasons why the original guidelines have had poor adherence and  
19 present an approach to improve their uptake. We present in this paper an R markdown  
20 reporting template file that guides the production of text and generates workflow diagrams  
21 based on user input. This R Markdown template contains, as an example in this instance, a  
22 set of minimum information requirements *specifically for the data pre-treatment and data*  
23 *analysis section* of biomarker discovery metabolomics studies, (gleaned directly from the  
24 original proposed guidelines by Goodacre et al). These minimum requirements are presented  
25 in the format of a questionnaire checklist in an R markdown template file. The R Markdown  
26 reporting template proposed here can be presented as a starting point to encourage the data  
27 analysis section of a metabolomics manuscript to have a more logical and stepwise  
28 presentation and to contain enough information to be understandable and reusable. The idea  
29 is that these guidelines would open to user feedback, modification and updating by the  
30 metabolomics community via GitHub.

31

32 **Keywords:** reproducibility; minimum guidelines; reporting; data analysis; reporting

33

34

#### 35 **1. Introduction**

36 Metabolomics data mining describes the application of strategic data analysis methods  
37 incorporating artificial intelligence, machine learning, statistics and database operations to  
38 extract meaningful and useful information from high-dimensional and high-volume

39 metabolomics datasets. This is a complex, time-consuming process including many steps  
40 with many possible options at each step. Metabolomics analysis is complicated by the  
41 metabolome's vast complexity and dynamics. Metabolites are present in a wide range of  
42 concentrations, from low abundance signaling molecules to high abundance compounds of  
43 central metabolism. Metabolites are also subject to temporal and spatial variability and are  
44 affected by environmental influences such as circadian and diet fluctuations to name a  
45 few. Further complexity is encountered in metabolomics investigations due to the fact that  
46 data analysis is often based on open source tools, each having their own parameter  
47 dependencies, also metabolomics datasets typically contain missing values and the handling  
48 of these can greatly influence the result of downstream analysis [1, 2]. For detailed  
49 discussion and reviews of the data analysis step in metabolomics and its complexities the  
50 reader is referred to the following publications [3-7].

51 Despite the obvious complexity and importance of the data mining step in the overall  
52 pipeline of any metabolomics study, this section of the workflow is often given scant  
53 attention in the write up of scientific research articles. Metabolomics data analysis sections  
54 have been found to be plagued by inconsistent reporting specifically with regards to  
55 structure, details reported, and performance metrics used [3].

56 No standard method for how to analyse metabolomics data exists and therefore data analysis  
57 is in constant evolution with new methods frequently being proposed in the literature. To  
58 discover the best methods, to build upon existing approaches and to conduct meta-analysis,  
59 the data analysis write up in metabolomics studies needs to be understandable and imitable,  
60 at a minimum. Furthermore, for those new to metabolomics data analysis, the starting point  
61 to construct a data analysis plan would most likely involve examining previously published  
62 research in the same field with a view to reusing or adapting the various approaches. For  
63 these purposes the current standard of reporting of the data analysis sections in metabolomics  
64 studies manuscripts is woefully insufficient. The immediate improvement of reporting of the  
65 data analysis step is therefore vital to advance understanding and to promote reuses of data  
66 analysis protocols and eventually move closer to the ideal of reproducibility.

67 Minimum reporting guidelines for data analysis in metabolomics [8] were first published in  
68 2007. These comprehensive guidelines cover: 1: Design of Experiment (sample  
69 collection/matching, and data acquisition scheduling of samples); 2: Data Collection 3: Data  
70 pre-processing (data cleaning, outlier detection, row/column scaling, or other  
71 transformations; Definition and parameterization of subsequent visualizations) 4: Data  
72 pre-treatment (row wise and column wise operations such as normalisation, scaling,  
73 centering and transformation to make data more amenable to statistical analysis); and finally,  
74 4: Actual data analysis which includes algorithm selection, univariate analysis and  
75 multivariate analysis. These reporting guidelines were published as an article but were not  
76 subsequently published in the format of a guidelines checklist with an explanation and  
77 elaboration document nor were they formally disseminated. Results of our recent review on  
78 reporting of the actual data analysis step in metabolomics indicate that these original  
79 reporting guidelines have had very poor take up, at least for the data pre-treatment and actual  
80 analysis section of metabolomics studies [2]. For example 89% (23 out of a total of 25) of

81 studies reviewed from the years 2008 to 2014 did not mention the proportion of missing  
82 values nor how the missing values were dealt with. Less than half the studies reviewed  
83 reported on any kind of quality control procedure and less than half had any mention of  
84 outlier detection and/or removal.

85 Reasons why those original reporting guidelines for metabolomics have had such poor take  
86 up may include the fact that they never progressed from the “proposed” stage to being  
87 formally published as practical guidelines along with a detailed “explanation and elaboration  
88 document” and they are not required by most journals for publishing a metabolomics  
89 manuscript. Also their comprehensiveness may have inhibited their uptake as there may have  
90 been an overwhelming amount of information to work through.

91 Strategies to increase the uptake and impact of guidelines can be adopted by their authors  
92 such as publishing the guidelines in multiple journals to ensure quicker and wider  
93 dissemination; also authors can approach journals and ask them to include the guidelines in  
94 their “Instructions to authors” section and publish commentaries to endorse them [9]. Official  
95 society and community wide recommendations can also positively influence the uptake of  
96 guidelines. There are a number of recommendations on a dedicated page of the EQUATOR  
97 (Enhancing the QUALity and Transparency Of health Research) Network website [10] on  
98 how to effectively disseminate your reporting guideline:  
99 [http://www.equator-network.org/toolkits/developing-a-reporting-guideline/disseminating-y-](http://www.equator-network.org/toolkits/developing-a-reporting-guideline/disseminating-our-reporting-guideline/)  
100 [our-reporting-guideline/.](http://www.equator-network.org/toolkits/developing-a-reporting-guideline/disseminating-our-reporting-guideline/)) This would of course be in addition to previous recommendations  
101 of supplementary audit, open code and script sharing [11].

102 In recognition of the complexity of data analysis in metabolomics we propose that distinct  
103 reporting guidelines be drawn up for separate sections of the data analysis pipeline (design of  
104 experiment, data reduction and deconvolution, data pre-processing, data annotation and  
105 identification, data pre-treatment and data analysis) as their various steps are often carried out  
106 at different times, (sometimes years pass between different steps), by different individuals or  
107 groups of various skill sets, or even often outsourced to different locations. Adherence to  
108 guidelines is more likely to be achieved by discretising each section of the pipeline into  
109 succinct guideline sets (modules) which can be adopted by the relevant analyst(s) *at the time*  
110 *of manuscript writing* and incorporated into the report.

111 In other established omics reporting guidelines’, namely MIAME [12] for microarrays and  
112 MIAPE [13] for proteomics, instructions do not exist for the data analysis part of the  
113 pipeline. The Equator (Enhancing the QUALity and Transparency Of health Research)  
114 Network [10] is an international initiative aimed at promoting transparent and accurate  
115 reporting of health research studies to enhance the value and reliability of medical research  
116 literature. The Equator Network does not contain any guidelines for reporting of multivariate  
117 data analysis/high dimensional data analysis/omics data analysis/supervised data analysis.  
118 The biosharing website MIBBI (Minimum Information for Biological and Biomedical  
119 Investigations) [14] does not contain any standards for data analysis reporting but it does  
120 reference a standard called CIMR- Core Information for Metabolomics Reporting (CIMR)  
121 [15] which refers back to the original proposed guidelines [8]. However since these initial

122 proposed guidelines in 2007 no further work has been published on the development or  
123 dissemination of metabolomics data analysis reporting guidelines.

124 Since the Metabolomics Standards Initiative (MSI) [16] significant developments in data  
125 reporting standards in metabolomics have been made through many initiatives including  
126 COSMOS [17], MetaboLights [18] and FAIR [19] which endeavour to ensure consistency  
127 of metadata between datasets, and facilitate data reuse and data merger across studies [20].  
128 However with regards to reporting of the data *analysis* of metabolomics studies, since the  
129 original guidelines [8] there have been no further advancements.

130 There has recently been a proliferation of reporting guidelines in biomedical research [21],  
131 there are currently 407 reporting guidelines on the Equator Network [10] many containing  
132 extensions and different versions. However despite this, compliance levels with these  
133 guidelines have been disappointing [22]. It has been noted that the “main problem”  
134 preventing the uptake of guidelines is that they are used too late in the research process, when  
135 it is too late to discover important things that have been missed or could have been done  
136 better [23].

## 137 1.2 Data analysis reporting using R Markdown

138

139 With this information in mind we suggest that adherence to guidelines could be facilitated if  
140 reporting guideline modules were contained in authoring tools such as the one we present in  
141 this study using R Markdown. R Markdown is a free and open source authoring framework  
142 [24]. R Markdown documents are fully reproducible and support dozens of static and  
143 dynamic output formats. A single R Markdown file can be used to both save and execute  
144 code and generate high quality reports that can be shared with an audience. R Markdown  
145 starts with a plain text file that is edited by the user that has the extension .Rmd. This plain  
146 text file then generates a new file that contains user selected text, code, and results from the  
147 .Rmd file. The new file can be a finished web page, PDF, MS Word document, slideshow,  
148 notebook, handout, book, dashboard, package vignette or other format.

149 As efforts towards computational reproducibility continue, an authoring tool in R markdown  
150 such as the one we present has the advantage that it can simultaneously achieve the aims of  
151 both ensuring reporting standards are adhered to while also having the potential to embed the  
152 code to perform the analysis. Such an authoring tool could therefore ultimately provide an  
153 uninterrupted and transparent workflow from the initial stage of data in to the final output of  
154 an analysis report. Our example of an authoring tool in this instance is solely focused on  
155 reporting of the statistical data analysis step of the pipeline, incorporating the data  
156 pre-treatment step and the actual data analysis step.

157

158

159

160

### 161 1.3 Objectives

- 162 1. To present a set of previously proposed minimum reporting guidelines in the form  
163 of a checklist specifically for the data analysis step of metabolomics biomarker  
164 discovery studies. There are typically 4 phases to this data analysis pipeline,  
165 although aside from pre-treatment the other steps are not essential but are  
166 commonly used
- 167 • data pre-treatment
  - 168 • univariate data analysis to identify significant features that  
169 differ between groups
  - 170 • unsupervised data analysis to discover correlated features or  
171 identify hidden subgroups or to visualise separation and  
172 identify outliers
  - 173 • Supervised data analysis, specifically for developing  
174 prediction models and/or biomarker identification.
  - 175 • Receiver Operating Curve (ROC) analysis
- 176
- 177 2. To provide an authoring tool to promote standardised comprehensive  
178 reporting on data analysis which will also generate workflow diagrams.

## 179 2. Methods

### 180 2.1: Checklist of minimum information for reporting data analysis in metabolomics

181 The development of a reporting guideline checklist, specifically for the data pre-treatment  
182 and data analysis sections of the metabolomics pipeline, with a view to general applicability  
183 to other omic domains.

184 This checklist was compiled based on the complete information required to construct a  
185 workflow diagram and to repeat the analysis using the reader's own version of code. The  
186 main areas of omissions which lead to confusion and ambiguity when conducting our review  
187 [3] helped to inform this checklist.

188 Existing guidelines which helped to shape this guideline list included the TRIPOD [25]  
189 statement, GRIPS [26] statement and REMARK guidelines [27]. Of course, the main  
190 document informing these guidelines is the original proposed guidelines for data analysis  
191 reporting in metabolomics by Goodacre et al [8] which covers the reporting of every part of  
192 the data analysis of a metabolomics experiment from design of experiment through  
193 pre-processing to data pre-treatment and final analysis.

### 194 2.2: An authoring tool for reporting statistical analysis of predictive omics

195 The development of an authoring tool using R Markdown. This reporting guideline checklist  
196 is presented as a questionnaire in an R markdown file that guides the production of text and  
197 workflow diagrams based on user input. These reporting guidelines are intended to form a  
198 neutral and malleable framework and have general applicability and interoperability across  
199 various omics domains. These can be extended as needed by different domains or studies but

200 would represent a minimum set of information to be supplied whenever predictive data  
201 mining in metabolomics is carried out. We purposely do not develop a “user friendly” web  
202 interface as the goal is for users to operate within the R Markdown environment.

203

## 204 3. Results

205 3.1: Minimum Information about a Data Analysis (MIDAS) checklist (Guidelines checklist specifically for the  
206 data analysis step)

207 Guidelines of two types generally exist: guidelines for reporting and guidelines for protocols.  
208 Since the area of data mining for metabolomics is still nascent we suggest that guidelines or  
209 limitations on methodology at this point would be premature as the optimal methods for  
210 extracting clinically useful biomarkers has clearly not been established. Therefore our  
211 guidelines pertain only to reporting.

212

### 213 MIDAS Guidelines Checklist

#### 214 Pre-treatment

- 215 • What are the dimensions of the dataset entering this phase of analysis?
- 216 • What percentage of the data is missing values?
- 217 • Is imputation (I) performed?
- 218 • If yes describe method
- 219 • Is normalisation (N) performed?
- 220 • If yes describe method
- 221 • Is transformation (T) performed?
- 222 • If yes describe method
- 223 • Is scaling (S) performed?
- 224 • If yes describe method
- 225 • Is filtering (F) applied to the dataset at this point?
- 226 • If yes describe method
- 227 • Is a QC (QC) method employed on the dataset?
- 228 • Please describe
- 229 • Outline the order of the pre-treatment steps performed on the dataset
  - 230 • E.g. I-> T-> S->N->F->QC
- 231 • Have the dimensions of the dataset changed from the outset of pre-treatment to
- 232 the end of pre-treatment?
- 233 • Provide details on the package or program used for this phase of the analysis
- 234 • If in house code is used provide it or a link to it and also the language the code is
- 235 written in.

236

#### 237 Univariate analysis

- 238 • What are the dimensions of the dataset entering this phase of analysis?
- 239 • Is univariate testing performed?
- 240 • If yes describe method
- 241 • Is a multiple testing correction employed with this method?
- 242 • If yes describe method



- 243 • Are other methods of univariate testing performed?  
244 • If yes describe methods  
245 • Are multiple testing correction employed with these methods?  
246 • If yes describe method  
247 • Please report p-values and adjusted p-values.  
248 • Please report test statistics and confidence intervals.  
249 • Have the dimensions of the dataset changed from the outset of univariate analysis  
250 to the end of univariate analysis? If yes provide the dimensions of the dataset at  
251 the end of univariate analysis and make clear how the dimensions have changed  
252 • Provide details on the package or program used for this phase of the analysis  
253 • If in house code is used provide it or a link to it and also the language the code is  
254 written in.  
255 • If a list of potential biomarkers is produced at this point please state this explicitly  
256

### 257 **Unsupervised analysis**

- 258 • What are the dimensions of the dataset entering this phase of the analysis?  
259 • Are unsupervised methods employed for visualisation and/ or data reduction  
260 and/or correlation analysis?  
261 • If yes describe the algorithm used.  
262 • Is outlier detection and removal addressed at this point? If yes please describe and  
263 specify the outliers removed.  
264 • Are unsupervised analysis methods used for clustering?  
265 • If yes describe and provide distance metric.  
266 • Have the dimensions of the dataset changed? If yes how and why?  
267 • Provide the dimensions of the dataset at the end of unsupervised analysis.  
268 • Provide details on the package or program used for this phase of the analysis  
269 • If in house code is used provide it or a link to it and also the language the code is  
270 written in.  
271 • If a list of potential biomarkers is produced at this point please state this explicitly  
272

### 273 **Supervised analysis**

- 274 • What are the dimensions of the dataset at this point?  
275 • Are supervised methods employed?  
276 • If yes describe the supervised analysis described fully enough to allow imitation  
277 of the exact procedure. This would require reporting all the following  
278 information: all parameters; details of how data is split; details of how internal  
279 validation is conducted (i.e. Cross Validation); details of how meta-parameter  
280 optimization is performed; details about the chosen metric for assessing the  
281 predictive ability of the model and finally the overall description of the workflow.  
282 • Is more than one supervised method employed?  
283 • If yes describe the implementation of the other algorithm(s) fully enough to allow  
284 imitation of the exact procedure. This would require reporting all the following



285 information: all parameters; details of how data is split; details of how internal  
286 validation is conducted (i.e. Cross Validation); details of how meta-parameter  
287 optimization is performed; details about the chosen metric for assessing the  
288 predictive ability of the model and finally the overall description of the workflow.  
289 • Is external validation employed?  
290 • If yes describe the source of external data. Is the data from the same location/  
291 lab/timeline or a hold-out set from the original data?  
292 • Provide a confusion matrix of results  
293 • Provide results as average of N leave-multiple-out and external predictions  
294 • Are potential biomarkers identified? If yes, list them.  
295 • Have the dimensions of the dataset changed? If yes how and why?  
296 • Provide the dimensions of the dataset at the end of supervised analysis.  
297 • Provide details on the package or program used for this phase of the analysis  
298 • If in house code is used provide it or a link to it and also the language the code is  
299 written in.  
300 • If a list of potential biomarkers is produced at this point please state this  
301 explicitly.  
302

### 303 **Receiver Operating Curve (ROC) Analysis**

304 • Is ROC analysis performed on the identified putative biomarkers?  
305 • If yes please report on AUC, sensitivity and specificity  
306 • Provide details on the package or program used for this phase of the analysis  
307 • If in house code is used provide it or a link to it and also the language the code is  
308 written in.

309 For data analysis methods currently not covered here (for example, cluster analysis and other  
310 classification and feature selection methods) similar templates can be generated with their  
311 required parameter reporting and add it to the existing templates via adding branches to the  
312 GitHub repository. We actively invite participation from metabolomics community users to  
313 become involved in this collaborative venture.

314

315 *3.2 Link to GitHub repository containing markdown template.*

316 *<https://github.com/MSI-Metabolomics-Standards-Initiative/MIDAS>*

317 This R Markdown file containing our authoring template is contained in a GitHub repository.  
318 Having started as a code developer's collaborative platform, GitHub [28] is now the largest  
319 online storage space of collaborative works that exists in the world which makes it the ideal  
320 platform to share this R Markdown template file.

321 The beauty of R Markdown as stated above is that it can embed and execute code and this  
322 code can then be hidden or displayed in the final document. Even in the most basic of report  
323 writing templates such as the one presented here this is very useful as we can use

324 DiagrammeR [29], a flexible and powerful R package for generating graph and flowchart  
325 diagrams. It is necessary to state that because these diagrams depend on HTML and  
326 JavaScript for rendering they can only be used in HTML based output formats (they don't  
327 work in PDFs or MS Word documents). We can get around this by saving our workflow  
328 diagrams within RStudio as an image (JPEG, PNG, BMP, etc.) and inserting them into the  
329 final PDF or Word document version of the report.

330 The authoring tool presented here does not dictate or control the report produced by the data  
331 analyst and the report produced can continue to be edited after the final report is generated. It  
332 merely serves as a guide for the writer to construct their analysis report by reminding them of  
333 points to include. To our knowledge this is the first instance that such a markdown tool has  
334 been proposed to aid and formalize reporting guideline uptake.

335 This authoring template is currently available on the Metabolomics Standards Initiative  
336 GitHub Repository and is open and welcoming to extension, modification and improvement  
337 from the metabolomics community and as such is considered a work in progress and a  
338 dynamic tool.

339

#### 340 4. Discussion

341 Currently a copy-and paste paradigm in which results are generated in a statistical package  
342 and copied and pasted to a report document dominates data analysis reporting. Eventually a  
343 complete move away from this antiquated copy-and-paste system which is error prone and  
344 enables selective reporting is needed in order to fulfil the requirements of reproducible  
345 research. However, in the meantime, in this instance, our version of the MIDAS reporting  
346 template allows users to manually input results that they have obtained using other software,  
347 whilst also having the potential to contain fully executable code. This is so as to not exclude  
348 users other than R users from benefiting from using this reporting template as an authoring  
349 tool and to encourage the first steps towards reproducible research.

350 Employing R Markdown to help uptake of minimum reporting standards goes further than  
351 providing a checklist. By encouraging scientists to consider reporting standards at the time of  
352 manuscript writing it actively helps authors adhere to guidelines. These guidelines should be  
353 viewed by authors as helpful to the writing process as opposed to a "yet another hurdle along  
354 the journey to publication" [21]. Also, requiring that parameter choices are revealed in  
355 reporting, even default ones from online data analysis tools, will encourage non experts to  
356 deliberate on the choices they are making regarding the appropriate algorithms and  
357 parameters for the type of analysis that they are doing, which will further the advancement of  
358 the field.

359 We purposely do not provide a web end user interface for accessing this R Markdown  
360 template as we believe that data analysts need to become comfortable in environments such  
361 as R Markdown if the production of reproducible research papers is to become a reality.

362 Furthermore we feel that anyone who is capable of data analysis is more than capable of  
363 using R Markdown without the need of a “user friendly” web interface.

364 A modularised system of authoring tools would encourage the uptake of guidelines on two  
365 fronts: The modularised facet would ensure that each researcher at each stage of the pipeline  
366 would be responsible for following the appropriate guidelines pertaining to *their own* area of  
367 expertise. The authoring tool part would ensure that guidelines are addressed *at the time of*  
368 *writing* up that particular section, as opposed to a set of rules to consider for application to the  
369 manuscript just before journal submission when the entire article has already been written.

370 Ideally, it is envisaged that such a modularised system of reporting guideline authoring tools  
371 would evolve in the omics community whereby these modules could be concatenated as  
372 needed depending on the experiment, each module corresponding to a stage in the workflow  
373 pipeline, with all modules being extensible and modifiable according to the domain and  
374 experiment in question.

375 In computational biology extensibility and modifiability of tools are essential so that new  
376 methods can develop and build on the old ones without repetition or reinventing the wheel.  
377 For this reason this R Markdown file is not presented here as an end result but is proposed as  
378 a starting point to encourage the data analysis section of metabolomics papers to have a more  
379 logical and stepwise presentation and to contain enough information to be understandable.  
380 So, even though this R Markdown file only attends to the authoring and not the analysis of  
381 metabolomics data we hope that it will coax data analysts into the environment of R  
382 Markdown (and GitHub) and therefore be a nudge along the road towards readable, and  
383 ultimately, reproducible, metabolomics research.

384 Here are the instructions to use this R Markdown authoring template.

- 385 1. Go to the **GitHub** repository:  
386 <https://github.com/MSI-Metabolomics-Standards-Initiative/MIDAS>
- 387 2. Click the “clone or download” button on the right hand side of the page and download  
388 the folder as a zip file.
- 389 3. Download latest version of **R Studio** if you do not have it.
- 390 4. Open the folder and open the *MIDAS.rmd* file in **R studio**.
- 391 5. Start editing and writing the report of your data analysis guided by the questions in  
392 green directly inside the *MIDAS.rmd* file.
- 393 6. After the pre-treatment, univariate analysis, unsupervised analysis and supervised  
394 analysis sections have been completed the next section is to produce workflow  
395 diagrams.
- 396 7. Follow the instructions in green to produce a workflow diagram of pre-treatment  
397 steps.
- 398 8. Click on the knit button and knit to HTML to see how the generated report looks.
- 399 9. Knit to PDF or Word to render the report to a pdf or word document as you wish.
- 400 10. PDF and Word reports will not contain the diagrams so these need to be saved in the  
401 viewer pane as an image (JPG /BMP etc) to your local folder

- 402 11. Insert the workflow diagrams into your Word or PDF report that you have saved to  
403 your local folder.
- 404 12. Render the document to HTML and workflow diagrams will be included anyway.

405 **Author Contributions:** For research articles with several authors, a short paragraph specifying their individual  
406 contributions must be provided. The following statements should be used “conceptualization, E.C;  
407 methodology, E.C.; software, E.C.; validation, RS; writing—original draft preparation, EC an.; writing—review  
408 and editing, R.S.”,

409 **Funding:** This research was funded by Science Foundation Ireland

## 410 References

- 411 1. Gromski, P.S., et al., *Influence of missing values substitutes on multivariate analysis*  
412 *of metabolomics data*. *Metabolites*, 2014. **4**(2): p. 433-452.
- 413 2. van den Berg, R.A., et al., *Centering, scaling, and transformations: improving the*  
414 *biological information content of metabolomics data*. *BMC Genomics*, 2006. **7**: p.  
415 142-142.
- 416 3. Considine, E.C., et al., *Critical review of reporting of the data analysis step in*  
417 *metabolomics*. *Metabolomics*, 2018. **14**(1): p. 7.
- 418 4. Cambiaghi, A., M. Ferrario, and M. Masseroli, *Analysis of metabolomic data: tools,*  
419 *current strategies and future challenges for omics data integration*. *Briefings in*  
420 *bioinformatics*, 2017. **18**(3): p. 498-510.
- 421 5. Bartel, J., J. Krumsiek, and F.J. Theis, *Statistical methods for the analysis of*  
422 *high-throughput metabolomics data*. *Computational and structural biotechnology*  
423 *journal*, 2013. **4**(5): p. e201301009.
- 424 6. Ren, S., et al., *Computational and statistical analysis of metabolomics data*.  
425 *Metabolomics*, 2015. **11**(6): p. 1492-1513.
- 426 7. Tugizimana, F., et al., *A conversation on data mining strategies in LC-MS*  
427 *untargeted metabolomics: Pre-processing and pre-treatment steps*. *Metabolites*,  
428 2016. **6**(4): p. 40.
- 429 8. Goodacre, R., et al., *Proposed minimum reporting standards for data analysis in*  
430 *metabolomics*. *Metabolomics*, 2007. **3**(3): p. 231-241.
- 431 9. Simera, I., et al., *Guidelines for Reporting Health Research: The EQUATOR*  
432 *Network's Survey of Guideline Authors*. *PLOS Medicine*, 2008. **5**(6): p. e139.
- 433 10. *The Equator Network*.
- 434 11. Meier, R., et al., *Bioinformatics can boost metabolomics research*. *Journal of*  
435 *biotechnology*, 2017. **261**: p. 137-141.
- 436 12. Brazma, A., et al., *Minimum information about a microarray experiment*  
437 *(MIAME)-toward standards for microarray data*. *Nat Genet*, 2001. **29**(4): p. 365-71.
- 438 13. Taylor, C.F., et al., *The minimum information about a proteomics experiment*  
439 *(MIAPE)*. *Nat Biotech*, 2007. **25**(8): p. 887-893.
- 440 14. *The Biosharing Website*. Available from:  
441 [https://biosharing.org/standards/?selected\\_facets=isMIBBI:true](https://biosharing.org/standards/?selected_facets=isMIBBI:true).
- 442 15. *CIMR-Core Information for Metabolomics Reporting*.
- 443 16. Fiehn, O., et al., *The metabolomics standards initiative (MSI)*. *Metabolomics*, 2007.  
444 **3**(3): p. 175-178.

- 445 17. Salek, R.M., et al., *COordination of Standards in MetabOlomicS (COSMOS):*  
446 *facilitating integrated metabolomics data access.* Metabolomics, 2015. **11**(6): p.  
447 1587-1597.
- 448 18. Steinbeck, C., et al., *MetaboLights: towards a new COSMOS of metabolomics data*  
449 *management.* Metabolomics, 2012. **8**(5): p. 757-760.
- 450 19. Wilkinson, M.D., et al., *The FAIR Guiding Principles for scientific data management*  
451 *and stewardship.* 2016. **3**: p. 160018.
- 452 20. Spicer, R.A., R. Salek, and C. Steinbeck, *Compliance with minimum information*  
453 *guidelines in public metabolomics repositories.* 2017. **4**: p. 170137.
- 454 21. The, P.M.E., *From Checklists to Tools: Lowering the Barrier to Better Research*  
455 *Reporting.* PLoS Medicine, 2015. **12**(11): p. e1001910.
- 456 22. Glasziou, P., et al., *Reducing waste from incomplete or unusable reports of*  
457 *biomedical research.* Lancet, 2014. **383**(9913): p. 267-76.
- 458 23. Marusic, A., *A tool to make reporting checklists work.* BMC Medicine, 2015. **13**(1):  
459 p. 243.
- 460 24. Baumer, B. and D. Udwin, *R Markdown.* Wiley Interdisciplinary Reviews:  
461 Computational Statistics, 2015. **7**(3): p. 167-177.
- 462 25. Moons, K.G., et al., *Transparent Reporting of a multivariable prediction model for*  
463 *Individual Prognosis or Diagnosis (TRIPOD): explanation and elaboration.* Annals of  
464 Internal Medicine, 2015. **162**(1): p. W1-73.
- 465 26. Janssens, A.C.J.W., et al., *Strengthening the Reporting of Genetic Risk Prediction*  
466 *Studies: The GRIPS Statement.* PLoS Medicine, 2011. **8**(3): p. e1000420.
- 467 27. McShane, L.M., et al., *REporting recommendations for tumour MARKer prognostic*  
468 *studies (REMARK).* British Journal of Cancer, 2005. **93**(4): p. 387-391.
- 469 28. *GitHub.* Available from: <https://github.com/>.
- 470 29. Sveidqvist, K., et al., *DiagrammeR: Create Graph Diagrams and Flowcharts Using R.*  
471 R package version 0.9. 0. URL: <https://CRAN.R-project.org/package=DiagrammeR>,  
472 2017.