

The Hitchhiker's Guide to Statistical Analysis of Feature-based Molecular Networks from Non-Targeted Metabolomics Data

Abzer K. Pakkir Shah^{1,2}, Axel Walter^{1,2,3}, Filip Ottosson⁴, Francesco Russo⁴, Marcelo Navarro-Diaz², Judith Boldt^{1,5}, Jarmo-Charles J. Kalinski^{1,6}, Eftychia E. Kontou^{1,7}, James Elofson⁸, Alexandros Polyzois^{1,9}, Carolina González-Marín^{1,10}, Shane Farrell^{11,12}, Marie R. Aggerbeck^{1,13}, Thapanee Pruksatrakul^{1,14}, Nathan Chan¹⁵, Yunshu Wang¹⁵, Magdalena Pöchlhacker^{1,16}, Corinna Brungs¹⁷, Beatriz Cámara¹⁸, Andrés Mauricio Caraballo-Rodríguez¹⁹, Andres Cumsille¹⁸, Fernanda de Oliveira^{20,21}, Kai Dührkopp²², Yasin El Abiead²³, Christian Geibel², Lana G. Graves^{24,25}, Martin Hansen¹³, Steffen Heuckeroth²⁶, Simon Knoblauch², Anastasiia Kostenko⁸, Mirte C. M. Kuijpers²⁷, Kevin Mildau^{1,28,29}, Stilianos Papadopoulos Lambidis², Paulo Wender Portal Gomes²³, Tilman Schramm^{2,30}, Karoline Steuer-Lodd^{2,30}, Paolo Stincone², Sibgha Tayyab², Giovanni Andrea Vitale², Berenike C. Wagner², Shipei Xing²³, Marquis T. Yazzie⁸, Simone Zuffa^{23,31}, Martinus de Kruijff³², Christine Beemelmans^{32,33}, Hannes P. Link², Christoph Mayer², Justin J.J. van der Hoof^{1,29,34}, Tito Damiani¹⁷, Tomáš Pluskal¹⁷, Pieter C. Dorrestein^{23,31}, Jan Stanstrup³⁵, Robin Schmid^{1,17}, Mingxun Wang^{1,15}, Allegra T. Aron^{1,8}, Madeleine Ernst^{4,#}, Daniel Petras^{1,2,30,#}

1. Virtual Multi-Omics Laboratory, The Internet

2. University of Tuebingen, Interfaculty Institute of Microbiology and Infection Medicine, Tübingen 72076, Germany

3. Applied Bioinformatics, Department of Computer Science, University of Tübingen

4. Section for Clinical Mass Spectrometry, Danish Center for Neonatal Screening, Department of Congenital Disorders, Statens Serum Institut, Artillerivej 5, DK-2300 Copenhagen S, Denmark

5. Leibniz-Institute DSMZ - German Collection of Microorganisms and Cell Cultures, Braunschweig, Germany

6. Department of Biochemistry and Microbiology, Rhodes University, 6140, Makhanda, South Africa

7. The Novo Nordisk Foundation for Biosustainability, Technical University of Denmark, Kemitorvet 220, 2800 Kongens Lyngby, Denmark

8. Department of Chemistry and Biochemistry, University of Denver, USA

9. Boyce Thompson Institute and Department of Chemistry and Chemical Biology, Cornell University, Ithaca, New York 14853, USA

10. Universidad EAFIT, Carrera 49, Cl. 7 Sur #50, Medellín, Antioquia, Colombia

11. Bigelow Laboratory for Ocean Sciences, East Boothbay, ME 04544, USA

12. School of Marine Sciences, Darling Marine Center, University of Maine, Walpole, ME 04573, USA

13. Department of Environmental Science, Aarhus University, Frederiksborgvej 399, 4000 Roskilde, Denmark

14. National Center for Genetic Engineering and Biotechnology (BIOTEC), National Science and Technology Development Agency (NSTDA), Thailand Science Park, Pathum Thani, 12120, Thailand

43 15. University of California Riverside, Department of Computer Science, 900 University Ave,
44 Riverside CA, USA
45 16. Department of Food Chemistry and Toxicology, University of Vienna, Vienna, Austria
46 17. Institute of Organic Chemistry and Biochemistry of the Czech Academy of Sciences, Prague,
47 Czech Republic
48 18. Laboratorio de Microbiología Molecular y Biotecnología Ambiental, Centro de Biotecnología
49 DAL, Universidad Técnica Federico Santa María, Valparaíso, Chile
50 19. Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San
51 Diego, 9500 Gilman Dr., San Diego, CA, 92093, USA
52 20. Department of Biotechnology, Engineering School of Lorena, University of São Paulo, Lorena,
53 São Paulo 12602-810, Brazil
54 21. Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San
55 Diego, La Jolla, California 92093, USA
56 22. University of Jena, Department of Bioinformatics, Jena, Germany
57 23. Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San
58 Diego, 9500 Gilman Dr., San Diego, CA, 92093, USA
59 24. University of Tuebingen, Department of Environmental Systems Analysis, Tübingen 72076,
60 Germany
61 25. Leibniz Institute of Freshwater Ecology and Inland Fisheries (IGB), Berlin, Germany
62 26. Institute of Inorganic and Analytical Chemistry, University of Münster, Münster, Germany
63 27. Department of Ecology, Behaviour and Evolution, University of California San Diego, 9500
64 Gilman Dr., San Diego, CA, 92093, USA
65 28. Department of Analytical Chemistry, University of Vienna, Austria
66 29. Bioinformatics Group, Wageningen University & Research, 6708 PB Wageningen, the
67 Netherlands
68 30. University of California Riverside, Department of Biochemistry, 900 University Ave, Riverside,
69 CA, USA
70 31. Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and
71 Pharmaceutical Sciences, University of California San Diego, 9500 Gilman Dr., San Diego, CA,
72 92093, USA
73 32. Helmholtz Institute for Pharmaceutical Research Saarland (HIPS), Helmholtz Centre for
74 Infection Research (HZI), Campus E8, 66123 Saarbrücken, Germany
75 33. Saarland University, Saarbrücken 66123, Germany
76 34. Department of Biochemistry, University of Johannesburg, Johannesburg 2006, South Africa
77 35. Department of Nutrition, Exercise and Sports, University of Copenhagen, 1958 Frederiksberg
78 C, Denmark
79
80 # Correspondence should be addressed to Madeleine Ernst (MAET@ssi.dk) or Daniel Petras
81 (daniel.petras@uni-tuebingen.de)
82



83
84 **Feature-Based Molecular Networking (FBMN) is a popular analysis approach for LC-**
85 **MS/MS-based non-targeted metabolomics data. While processing LC-MS/MS data through**
86 **FBMN is fairly streamlined, downstream data handling and statistical interrogation is often**
87 **a key bottleneck. Especially, users new to statistical analysis struggle to effectively handle**
88 **and analyze complex data matrices. In this protocol, we provide a comprehensive guide**
89 **for the statistical analysis of FBMN results. We explain the data structure and principles of**
90 **data clean-up and normalization, as well as uni- and multivariate statistical analysis of**
91 **FBMN results. We provide explanations and code in two scripting languages (R and**
92 **Python) as well as the QIIME2 framework for all protocol steps, from data clean-up to**
93 **statistical analysis. Additionally, the protocol is accompanied by a web application with a**
94 **graphical user interface (<https://fbmn-statsguide.gnps2.org/>), to lower the barrier of entry**
95 **for new users. Together, the protocol, code, and web app provide a complete guide and**
96 **toolbox for FBMN data integration, clean-up, and advanced statistical analysis, enabling**
97 **new users to uncover molecular insights from their non-targeted metabolomics data. Our**
98 **protocol is tailored for the seamless analysis of FBMN results from Global Natural**
99 **Products Social Molecular Networking (GNPS and GNPS2) and can be adapted to other MS**
100 **feature detection, annotation, and networking tools.**

101

102 1. Introduction

103 Metabolomics aims to characterize and quantify the detectable spectrum of small molecules in
104 order to catalog and understand the metabolic dynamics within biological systems^{1,2}. Phenotypes
105 or environmental factors that distinguish samples within a given set are often reflected in the
106 chemical profiles of such small molecules across samples. Therefore, the characterization of
107 chemical distinctions and gradients between samples provides crucial information for describing
108 and understanding molecular mechanisms^{3,4}. Metabolomics studies usually employ targeted or
109 non-targeted approaches². Targeted metabolomics is typically hypothesis-driven and aims to
110 quantify known metabolites, often using internal standards and experimental methodology
111 optimized for the study. Non-targeted metabolomics, on the other hand, aims to detect a
112 maximum number of metabolites in order to comprehensively characterize the chemical profiles
113 within a given sample set.

114 To uncover molecular insights from non-targeted liquid chromatography tandem mass
115 spectrometry (LC-MS/MS) data, several software tools are available that assist with mining and
116 annotating metabolites, including feature detection and annotation tools⁵. Feature-Based
117 Molecular Networking (FBMN) is a popular analysis platform that integrates various feature-
118 detection tools with molecular networking for metabolite annotation and annotation propagation⁶
119 in the GNPS cloud ecosystem⁷. FBMN is routinely applied many biological disciplines, including
120 clinical studies^{8,9}, plant^{10–12} and environmental science^{13–16}, as well as microbiome^{17–19} and the
121 functional analysis of natural products^{20–22}. While platforms such as GNPS have improved the
122 way that we identify and characterize metabolites, the subsequent step — statistical analysis —
123 remains a challenge for many researchers. While resources like MetaboAnalyst^{23,24} provide
124 powerful solutions for the statistical analysis of metabolomics data, the complex multi-layer
125 information from FBMN results and other downstream annotation tools (e.g., SIRIUS) require
126 typically multiple matrix operations, data clean-up, normalization, before uni- and multivariate
127 statistical analyses. Most tools and analysis approaches that can be used to archive this are
128 typically custom scripts or different software tools that are scattered across different platforms.
129 This makes it especially difficult for users new to the field to effectively manage and analyze their
130 data. Moreover, while there are several tools available for individual clean-up and analysis steps
131 (see alternative approaches section), there is a lack of a comprehensive, user-friendly guidance
132 that covers the entire spectrum of data preparation and statistical analysis of FBMN results.
133 In this protocol, we provide a detailed guide that starts with FBMN results, offering an end-to-end
134 pipeline from feature detection, spectrum annotation, subsequent data clean-up and statistical
135 analysis steps. This step-by-step guide is complemented with ready-made code for the popular
136 statistical scripting and platforms R and Python, the QIIME2 toolkit ([https://github.com/Functional-](https://github.com/Functional-Metabolomics-Lab/FBMN-STATS)
137 [Metabolomics-Lab/FBMN-STATS](https://github.com/Functional-Metabolomics-Lab/FBMN-STATS)), as well as a web application designed to simplify the process
138 (<https://fbmn-statsguide.gnps2.org/>). The protocol provides thereby a seamless analysis guide for
139 FBMN results from the GNPS (<https://gnps.ucsd.edu>), and GNPS2 (<https://gnps2.org>) web
140 platforms, which can also be adapted to other MS feature detection and annotation tools.

141 **Feature-based Molecular Networking from LC-MS/MS Data**

142 Liquid chromatography-mass spectrometry (LC-MS) is one of the most prominent metabolomics
143 techniques, with applications in numerous research fields^{25–28}. Specifically, LC coupled with
144 tandem mass spectrometry (LC-MS/MS) has been widely used because it provides a broad
145 coverage of chemical space allowing for the simultaneous semi-quantitative detection and
146 qualitative annotation of many metabolites over a wide dynamic range^{14,29–32}. In addition to
147 providing the molecular mass, retention time, and isotopic pattern of a metabolite, MS/MS
148 provides structural information about the detected species. This is achieved through the
149 fragmentation of precursor ions into product ions and the measurement of their mass-to-charge
150 ratios (m/z) and abundances. This is usually done through Data-Dependent Acquisition (DDA),
151 where ions that are observed in MS1 survey scans are iteratively selected for further
152 fragmentation in subsequent MS/MS scans (See **Figure 1.1**). DDA operates by selecting the “top
153 N” peaks in each duty cycle, where “N” is a user-defined number. These peaks are chosen based
154 on their intensity and other user defined criteria through an automated process¹⁵. The resulting
155 MS/MS spectra of product ions can be used in several ways to determine a candidate structure:

156 1) via spectral matching against spectral libraries of experimental reference spectra or *in silico*
157 generated spectra^{33,34}; 2) via machine learning-based structural predictions using experimental
158 MS/MS-generated molecular fingerprints against structural databases^{35,36}; 3) and via *de novo*
159 structure prediction using molecular structure fingerprint prediction combined with neural
160 networks³⁷.

161 Non-targeted LC-MS/MS metabolomics is a powerful and versatile research approach that
162 enables high-throughput analysis and simultaneous detection of many small molecules, making
163 it an excellent method for gaining insights into biological systems (For more information on
164 Experimental Design and LC-MS/MS Data Acquisition, refer to **Box 1**). However, mining the vast
165 amount of data created by non-targeted metabolomics experiments remains a challenging task
166 despite a range of available resources that guide in the qualitative and quantitative aspects on
167 non-targeted metabolomics. Qualitative data exploration has been democratized by platforms
168 such as GNPS³⁸, by providing MS reference libraries, data analysis workflows, and compute
169 resources for the community. Molecular networking (MN) is GNPS' core concept and is based on
170 the comparison of all MS/MS spectra within a dataset by modification-aware similarity metrics,
171 which network features by their similar fragmentation patterns that are often reflective of structural
172 similarity. FBMN⁶ and Ion Identity Molecular Networking (IIMN)³⁹ add feature detection, improving
173 the (semi) quantitative quality within MN results. FBMN builds upon the classical MN by
174 harnessing both MS1 information, such as isotope patterns and retention time, and ion mobility
175 separation when used. FBMN can distinguish isomers with similar MS/MS spectra, which might
176 remain obscured in classical MN, through chromatographic or ion mobility separation.

177 IIMN enhances MS/MS-based spectral networks by adding connectivity based on the MS1 feature
178 shape correlation. It efficiently tackles the issue of unconnected ion adducts in Molecular
179 Networking by connecting ions from the same molecules into groups called ion identity networks
180 (IIN). This helps remove redundancy in MS-based metabolomics.

181

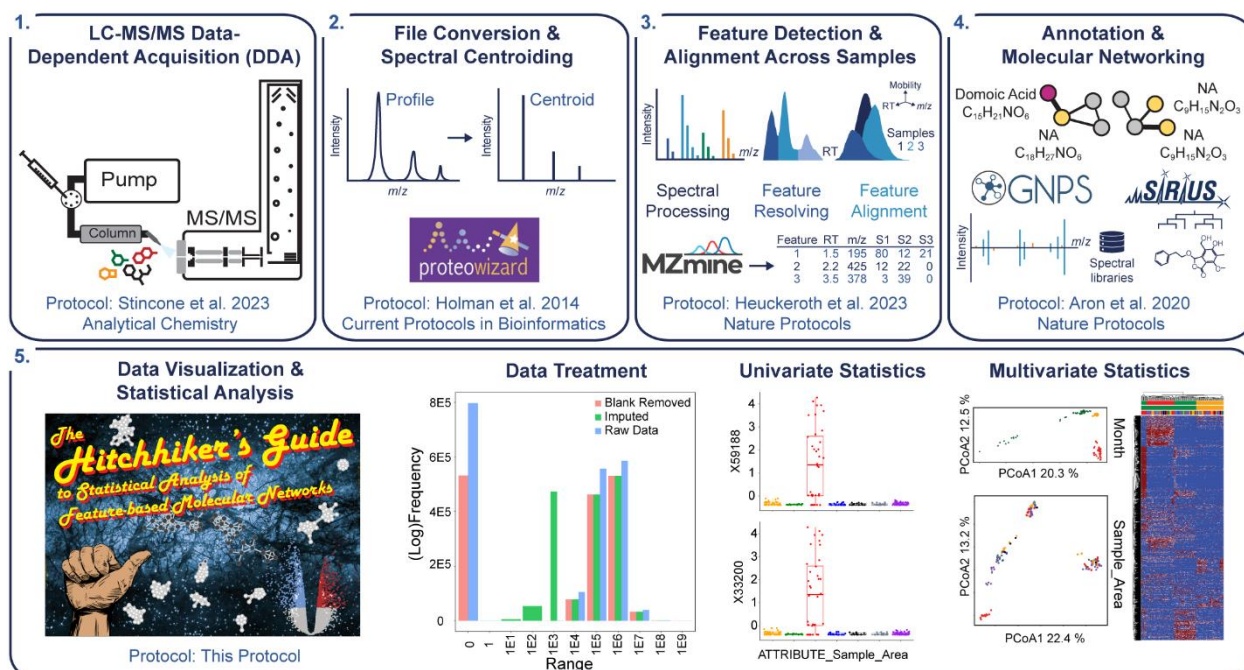
182 **Box 1 - Experimental Design and LC-MS/MS Data Acquisition**

To obtain high-quality and representative LC-MS/MS data, proper planning of the sampling and mass spectrometry analysis is essential. While this protocol article focuses on data analysis, we stress the importance of addressing the following topics prior to collecting any samples. It is crucial for researchers new to these experiments to seek guidance from a statistician and analytical chemists to guarantee optimal experimental design, instrument performance (e.g., system suitability tests), and analysis pipeline. Before proceeding with further processing and statistical analysis, raw LC-MS/MS data should be inspected by the user⁴⁰⁻⁴². For raw data processing, we recommend the MZmine protocol⁴³. Below, we provide a short checklist for guidance:

- **Experimental design and power calculation** are crucial when determining the suitable number of samples and replicates. In non-targeted metabolomics experiments, it is often challenging to predict certain values, like the feature **coefficient of variation** and the **expected effect size**, which are crucial for estimating the required sample size and the power of the study. To navigate this, reviewing previous studies with similar biological systems and research questions can provide an approximate estimation of these values.

As a general rule of thumb, when the effect size is smaller, one might need more samples or replicates.

- **Replicates (technical and biological)** to measure instrument and biological variance.
- **System and Process Blanks** to identify and correct for contamination that may be introduced during the sample collection, preparation, or measurement process. Some common blank samples include⁴⁴:
 - **Solvent blank:** Consists of only the solvent used to dissolve the sample. It is used to identify the contaminants present in the solvent. Also, adding this blank periodically in an analytical run reduces carryover. Blanks should be added into the same well plates or vials to cover similar contaminations.
 - **Extraction blank:** It is prepared by adding a known volume of solvent to a blank matrix such as water and extracting it the same way as a sample. This extracted blank is measured along the real sample to find the contaminants introduced during the extraction process.
- **Control samples**, e.g., negative and positive controls. Depending on the experimental design, control samples are essential and should be included in the same number as the treatment samples.
- **Quality Control (QC) samples** are needed to measure instrument performance. These can be in the form of pooled QC (for example, a combination of aliquots from each sample) or standard mixtures (such as a combination of reference standard chemical compounds or isotopically labeled compounds that can also serve as internal standards). These mixtures should span a broad chemical spectrum and cover a wide retention time range.
- **Randomization of sample injection order.** It is suggested to randomize the injection order throughout the samples. However, we recommend injecting blanks at the start of the queue to prevent carry-over, which could lead to the removal of actual features from the samples during the blank removal step. Depending on the experimental design, it might also be useful to select certain sample types with the injection order, e.g., KO (knockout) vs. WT (wild type) mutant strains or low vs. high biomass samples, to avoid carryover between them.
- **Internal Standards (IS)** can be added to every sample to track instrument performance, and if desired, quantify predefined metabolites. If no internal standards are available, “housekeeping features” such as ubiquitous contaminants or metabolites can be used to control for mass and retention time drift.



184
 185 **Figure 1: Flowchart of LC-MS/MS-based metabolomics experiment.** 1. Data-dependent
 186 acquisition of MS/MS spectra. 2. Centroiding and file conversion. 3. Feature detection. 4. Feature
 187 annotation, network propagation and clustering. 5. Data clean-up, statistical analysis, and
 188 visualization (blank removal, imputation, normalization and scaling, followed by data visualization
 189 and statistical analysis.)

190 Feature-based Molecular Networking Workflow

191 As highlighted in **Figure 1**, non/targeted LC-MS/MS analysis workflows typically consist of data
 192 acquisition, centroiding and file conversion, feature detection, and feature annotation, including
 193 spectrum/library matching, in silico spectrum annotation and annotation propagation and
 194 clustering approaches, such as spectral networking. From there, the resulting feature tables
 195 contain all metabolites / small molecules features detected, including quantitative information
 196 (e.g., peak area) in each sample measured. This resulting feature-sample matrix is then further
 197 processed by blank removal, imputation, normalization and scaling and finally data visualization
 198 and statistical analysis.

200 1. File Conversion

201
 202 Raw data acquisition in MS instruments entails first generating spectra in profile mode, also called
 203 continuous mode. In high-resolution instruments, each chemical entity is represented by signals
 204 of m/z ratios within a 5-20 ppm window, depending on the instrument's accuracy. The resulting
 205 peak profile typically approximates a Gaussian shape and is continuous. To reduce data
 206 complexity for downstream analysis, the data is simplified such that each peak in m/z dimension
 207 is represented by a single peak in the mass spectrum. This process is called centroiding or
 208 sometimes confusingly referred to as peak-picking, not to be confused with peak-picking in the
 209 chromatographic dimension termed "feature detection" below. Centroiding can either be

210 performed on-the-fly during acquisition by the vendor software or during file conversion using
211 tools like Proteowizard's msConvert⁴⁵ or ThermoRawFileParser⁴⁶ when converting from vendor-
212 specific formats into more accessible, community-driven formats such as mzML. When using the
213 "vendor" option in msConvert, the centroiding algorithm provided by the instrument's vendor is
214 used. These instrument-specific algorithms are typically more accurate than the algorithm
215 otherwise available in msConvert and are thus highly recommended (See **Figure 1.2**).

216

217 **2. Feature Detection**

218

219 The process of converting raw data into a table of putative metabolic features, along with their
220 relative abundances per sample, involves a pre-processing workflow that uses a series of
221 algorithms. The resulting table as shown in **Table 1** is referred to as a '**feature quantification**
222 **table**'. Open-source tools such as the R package XCMS⁴⁷ (often used with the package
223 CAMERA⁴⁸ for feature grouping), MZmine 3⁴⁹, MS-DIAL⁵⁰, or OpenMS⁵¹, in addition to vendor
224 specific tools can be utilized for this purpose.

225 For the present protocol, we decided to focus on MZmine 3. Firstly, it provides an interactive and
226 user-friendly graphical user interface (GUI) that can assist researchers without programming
227 skills. Secondly, the direct interfacing of MZmine 3 with the downstream annotation tools is
228 enabled by harmonized data exchange formats. Finally, the software offers a wizard for the
229 simplified generation of data-processing workflows, which reduces the number of parameters to
230 set and optimize for new users and experts. The process in MZmine 3 starts with importing
231 (centroided) mzML data (described in 'File Conversion') followed by the assembly of 'mass lists'
232 - i.e., lists of m/z values that exceed a user-defined intensity threshold. The workflow then
233 progresses through three main stages: feature detection, feature alignment, and feature
234 refinement (See **Figure 1.3**). For advanced optimization and fine-tuning, multiple tools such as
235 NeatMS⁵², MetaClean⁵³, and mzRAPP⁵⁴ exist to assess feature quality.

236

- 237 1. The feature detection phase is initiated by chromatogram building through the construction
238 of extracted ion chromatograms (EICs) by linking MS1 signals in consecutive scans based
239 on a maximum scan-to-scan mass deviation. This results in a list of features, each
240 characterized by a retention time (RT) and m/z value. Optional smoothing in the RT
241 dimension can be applied in the case of noisy data. The next step, 'Feature Resolving',
242 distinguishes between overlapping and co-eluting chromatographic peaks and is used to
243 link MS/MS spectra to their respective MS1 features. To remove redundant features
244 originating from isotopologues of the same parent ion, the ¹³C isotope filter can be
245 utilized. The 'Isotope Pattern Finder' identifies isotope signals of selected chemical
246 elements in each feature's mass list. The steps described above are carried out for the
247 feature list of each data file (sample) individually.
- 248 2. Next, in the feature alignment step, the individual feature lists created from multiple data
249 files are merged by matching features across all samples based on their RT and m/z
250 values. The peak alignment parameters are determined by the user and may differ
251 depending on the particular instrument used⁵⁵.
- 252 3. Lastly, the feature refinement phase can include a gap-filling procedure that accounts for
253 missing features in certain samples (e.g., signal below peak detection thresholds defined

254 in step 1). This procedure distinguishes between genuine absences and artifacts from the
255 feature-detection process. It locates signals in the original mzML (centroided) data by re-
256 evaluating their presence in individual samples for all the features in the merged feature
257 list and then replaces missing values with newly detected ones based on the user-defined
258 parameters such as *m/z* tolerance and RT tolerance. These tolerances set the window
259 within which the algorithm finds the new feature.

260 Furthermore, the merged feature list can be filtered by removing duplicate features,
261 features without a linked MS/MS spectrum, or features occurring outside a specified
262 retention time range (e.g., during re-equilibration phases). The final step involves
263 exporting the feature list as a feature quantification table (.csv). This table includes RT,
264 *m/z*, and relative abundance per sample for each feature. Additionally, a text file (.mgf) is
265 exported, describing the MS/MS and/or MS1 spectra linked to each feature. These output
266 files provide appropriate inputs for statistical analyses and data analysis in tools such as
267 GNPS and SIRIUS⁵⁶.

269 **Table 1. Example Feature Quantification Table.** The feature quantification table can be easily
270 converted to a table in a text format (example of a table of features).

	<i>m/z</i>	RT (min)	Adduct	Charge	Sample 1	Sample 2	...	Sample N
Feature 1	97.1082	4.6	[M+H] ⁺	+1	2.08e07	9.47e06	...	3.27e08
Feature 2	518.3032	2.0	[M+H] ⁺²	+2	1.88e07	5.56e05	...	2.11e06
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
Feature K	83.1017	1.6	[M+Na] ⁺	+1	4.77e04	8.13e03	...	5.17e09

271
272
273
274

3. Feature Annotation, Spectral Networking and Annotation Propagation

275 The Metabolomics Standard Initiative (MSI) outlines four levels of metabolite identification through
276 mass spectrometry to guide researchers in differentiating the level of identification rigor for the
277 reported metabolites⁵⁷. Level 1 annotations regard fully characterized compounds. To ensure
278 accurate annotation confidence at this level, it is necessary to have at least two independent
279 orthogonal data dimensions that match those of a pure compound analyzed in the same way
280 (e.g., precursor *m/z*, RT, MS/MS fragmentation pattern). There should be no contradictions in any
281 of the available data dimensions.

282 Level 2 confidence is assigned when data are matched against public or commercial data
283 libraries, such as MS/MS spectral matching on GNPS. Level 3 refers to compounds whose
284 chemical class can be putatively inferred through physicochemical features or data similarity with
285 known compounds (e.g., by spectral similarity (networking)), or using structural prediction tools

286 such as CSI:FingerID or CANOPUS^{35,37}. Finally, level 4 refers to unknown features that can be
287 consistently detected (e.g., defined *m/z* value, RT and MS/MS spectrum), but could not be
288 annotated through previous levels.

289 Feature annotation is essential in mass spectrometry-based metabolomics studies, especially to
290 understand the biological significance of the detected features. Feature annotation entails several
291 approaches, including database searches, spectral matching, and *in silico* annotation strategies.
292 *In silico* annotation tools, such as SIRIUS, MS2Query, Network Annotation Propagation (NAP),
293 Dereplicator, and Dereplicator+, predict metabolite identities based on spectral similarities⁵⁶ and
294 can only lead to MSI levels 2, 3 or 4.

295 Another innovative method that combines feature annotation with visualization is molecular
296 networking (MN), as shown in **Figure 1.4**. MN elucidates the structural relationships between
297 metabolites, highlighting potential biological pathways and processes. The utility of MN spans
298 across various fields, such as natural products, agriculture and clinical microbiology^{58–60}. Using
299 the MS-Cluster algorithm on the GNPS (<https://gnps.ucsd.edu>), and GNPS2 (<https://gnps2.org>)
300 web platforms, MN creates a molecular network by comparing spectral similarities between each
301 MS/MS spectra pair. With the same algorithm, GNPS allows the dereplication of MS/MS spectra
302 by comparing them against comprehensive spectral databases, enabling feature annotations of
303 varying reliability⁷. The .mzML or .mzXML spectra files can be analyzed on GNPS using the
304 classical MN workflow³⁸.

305 For more precise quantitative insights, FBMN has emerged as a significant advancement by
306 incorporating MS1 peak intensities, isotope patterns, retention times and ion mobility separation.
307 Consequently, FBMN distinguishes between isomers with near-identical fragmentation spectra,
308 but different retention times⁶. Unlike the classical MN, which required users to separately execute
309 molecular networking and MS1 analysis, FBMN conveniently accepts the output of feature
310 detection and alignment tools such as MZmine^{49,61} (see 'Feature Detection' above), MS-DIAL⁵⁰,
311 and XCMS⁴⁷, and is available on the GNPS web platform. This compatibility with other tools
312 makes FBMN seamlessly integrated into the overall analysis pipeline.

313

314 **4. Data Visualization and Statistical Analysis**

315 The feature quantification table (see **section 'Feature Detection'** and **Table 1**), contains a list of
316 features, such as *m/z* and RT, as well as their relative abundances per sample. This table
317 represents the basis dataframe for statistical analyses, which can help reveal distribution patterns
318 between sample types and determine which features are responsible for distinguishing between
319 them. The challenge lies in prioritizing the important features in a large dataset, considering
320 chemical and biological relevance, as well as statistical significance. An unsupervised approach
321 for initial exploration and visualization of the data is through dimensionality reduction techniques,
322 such as Principal Coordinates Analysis (PCoA). Ideally, such an approach will provide a 2- or 3-
323 dimensional plot, where similar samples are grouped together, apart from dissimilar samples.

324 Another unsupervised statistical approach is the use of hierarchical clustering to group samples
325 with similar relative abundance profiles of features. The results of such analysis are often
326 visualized in combination with a heatmap. This approach displays the features within each sample
327 colored according to their relative abundance, and groups them according to their similarity. A

328 dendrogram is drawn beside the heatmap to illustrate the hierarchical relationship between the
329 samples and features. Compound class ontologies such as ClassyFire⁶² or NPC⁶³ categorize
330 compounds based on shared structural features or biosynthetic origins and serve as high-level
331 annotations of the data. CANOPUS⁶⁴ predicts these compound classes from tandem mass
332 spectrometry data without searching in structure databases. Analyzing the distribution and variety
333 of compound classes, along with their up- and down-regulation, can yield biological insights that
334 may not be attainable when solely considering *m/z* and retention time values.

335 **Aim of the Protocol**

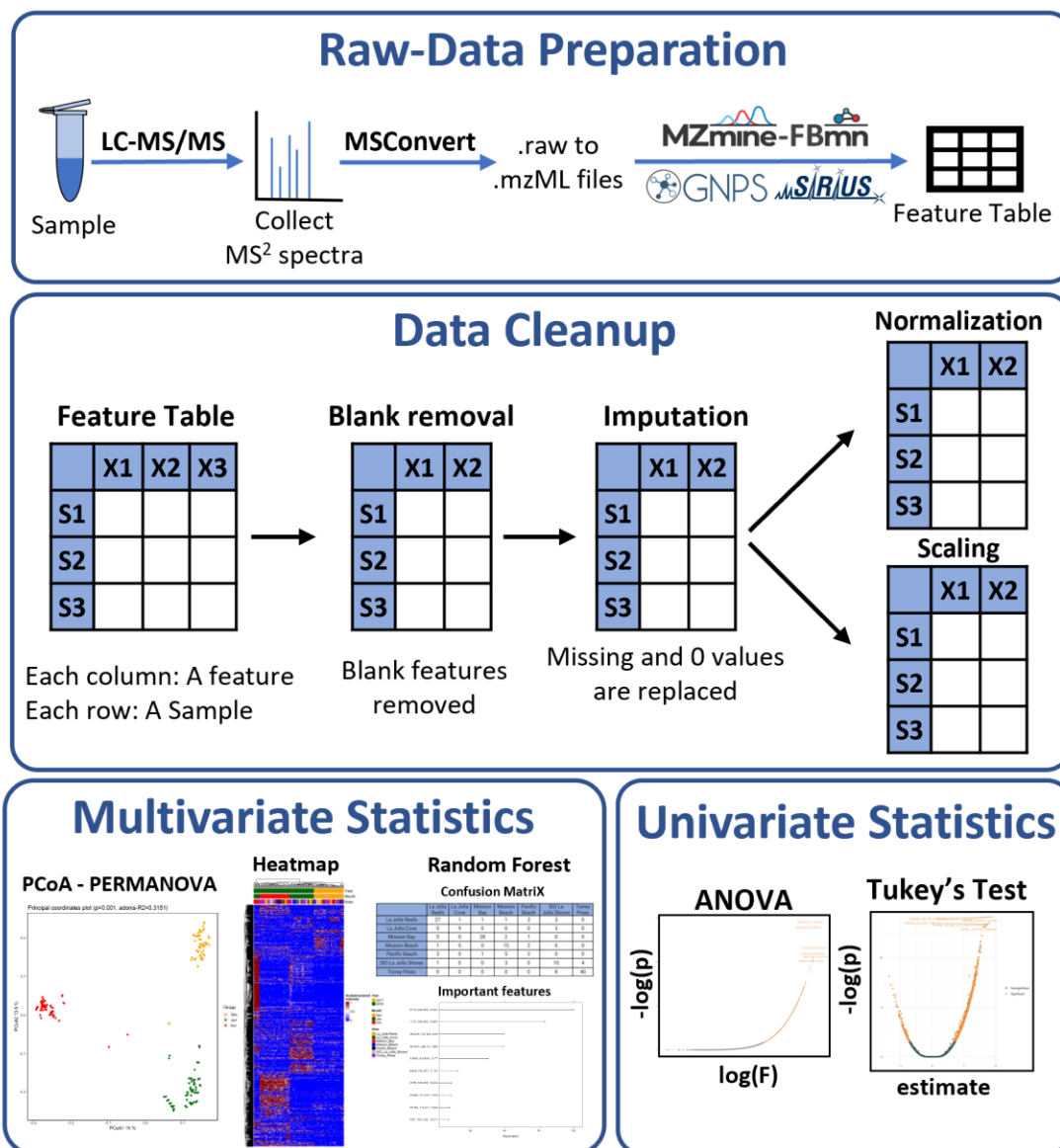
336 The goal of this protocol is to provide an integrated pathway for downstream data clean-up and
337 statistical analysis of FBMN results derived from non-targeted LC-MS/MS data (see **Figure 1.5**).
338 Integrating FBMN results with statistical analyses has poses several challenges, often
339 necessitating users to reformat, upload, and process the feature table with different external tools,
340 in order to ultimately manually combine the outcomes. Our approach addresses this gap by
341 offering a detailed guide and comprehensive solution to directly process and analysis the data
342 after FBMN in one pipeline, shown in **Figure 2**. This pipeline is provided in popular scripting
343 languages, R and Python, in conjunction with the well-known QIIME2 framework. It is available in
344 the form of Jupyter notebooks for local use and Google Colab notebooks for cloud-based
345 applications. Additionally, we have developed a web application with a graphical user interface
346 (GUI), which can be accessed at <https://fbmn-statsguide.gnps2.org/>. The main manuscript
347 focuses on the concepts and step-to-step guide for the R workflow, while the Supplementary
348 Information (SI) contains step-to-step guides for the Python, QIIME2, and Webapp workflows.
349 Though most steps are consistent across the workflows, any differences are addressed and
350 complemented with alternative solutions in the SI. This protocol is made for both newcomers and
351 experienced researchers in the metabolomics field:

- 352 ● **For Beginners:** It introduces essential tools, resources, and workflows. The guidelines
353 and code provided make it easier to understand common data processing and analysis
354 steps, facilitating navigation through the complexities of the field. The provided tools utilize
355 common programming languages (R, Python), the QIIME2 platform, and a GUI, allowing
356 users with diverse computational backgrounds to perform data analyses.
- 357 ● **For Experts:** It accelerates data analysis, ensuring faster interpretation of FBMN results.

358
359 As inputs, the protocol requires a feature table and its corresponding metadata table. Throughout
360 its execution, users receive:

- 361 ● Intermediate tables after each data cleanup step, aiding in comparison with the original
362 feature table. Tabular outputs for clustering results from Hierarchical Cluster Analysis
363 (HCA), list of statistically significant features as determined by ANOVA or Kruskal-Wallis
364 tests, and Random Forest outputs indicating feature importance. Significant features refer
365 to those that differ notably in at least one group when comparing multiple groups. Such
366 features can be further investigated to determine if they really cause the differences we
367 observed between groups or samples.
- 368 ● Visual outputs, such as PCoA score plots, heatmaps, volcano plots for significance tests,
369 and box plots, showcasing group differences for significant features.

370 This protocol helps with mapping the results of some of the statistical approaches (e.g., clustering,
 371 significant features) back to the FBMN network view. This is facilitated by importing these output
 372 feature tables, with feature IDs and the relevant information, into CytoScape in order to examine
 373 the molecular network. Moreover, as all our resources are publicly available on GitHub, users can
 374 actively raise issues or provide suggestions on GitHub.
 375



376
 377 **Figure 2: Overview of the Data Analysis Pipeline:** Integrating four core segments, the flowchart starts
 378 with sample collection and LC-MS/MS data acquisition, transitions to raw data conversion into mzML
 379 format, and results in generating a feature quantification table under “Raw Data Preparation”. This is
 380 followed by the “Data Cleanup” phase, emphasizing feature quantification table refinement, blank removal,
 381 imputation, and normalization strategies like Total-Ion-Current (TIC), Probabilistic Quotient Normalization
 382 (PQN) and scaling. Subsequently, the “Multivariate Statistics” segment showcases techniques such as
 383 PCoA plots, and heatmaps for effective data portrayal. In addition, the users are introduced to robust
 384 techniques including Random Forest classification. In the “Univariate Statistics” segment, tests such as
 385 ANOVA and Kruskal-Wallis test are discussed.

386 **Limitations and Challenges**

387 Our protocol for FBMN is aimed to offer advanced statistical analysis solutions for broad range of
388 users. We thus offer the notebooks and code in different scripting languages (R, Python, and
389 QIIME2) and platforms (Jupyter and Colab) as well as a web app to suit the specific needs and
390 preferences within the metabolomic community.

391 This broad range of choices, while useful, comes with its own set of challenges. For instance, in
392 the R Google Colab notebook, package installation can be time-consuming. Also, the inclusion of
393 readline commands, although beneficial for customization, can appear cryptic to beginners. On
394 the other hand, installing packages in the Python Google Colab notebook is relatively faster.
395 There is one vital point to note regarding the 'scikit-bio' package's incompatibility with Windows.
396 Thus, Windows OS users are advised to either use the Google Colab version or consider the
397 Windows Subsystem for Linux (WSL) for local operations.

398 Furthermore, while Google Colab stands as a user-friendly platform, it is not devoid of limitations.
399 One of the main concerns is that runtime automatically disconnects if the user leaves the Colab
400 session inactive for 90 minutes or after 12 continuous hours of usage. This leads to the loss of
401 the data and files they were working on from the Cloud session. Additionally, users must be aware
402 of the 77 GB disk space limitation and ensure timely downloading of their results.

403 Both the R and Python notebooks comprise over 70 steps, with a significant portion dedicated to
404 data organization. While these notebooks function smoothly with smaller datasets when run on
405 the Cloud, their performance can lag with larger datasets (e.g., those with over 100 samples and
406 more than 2,000 features), especially given the constraints of Google Colab. In such scenarios,
407 local execution is advisable. For local execution, we have provided guidance on using the
408 Anaconda Navigator, a user-friendly GUI platform, to set up Jupyter notebooks. However, MacOS
409 users might encounter installation challenges. As an alternative, MacOS users can opt for the 'pip
410 install' method. While numerous online resources can help with this, we have chosen not to delve
411 into the details here. The Streamlit WebApp for the protocol, although user-friendly, has its own
412 set of challenges. Notably, there's a data restriction of 200 MB, and larger datasets might
413 inadvertently slow down the app or even lead to server crashes.

414 Lastly, the QIIME2 notebook is broadly used and applicable for both the microbiome and
415 metabolomics communities. Our additional Jupyter notebook lets users import data directly from
416 a GNPS job link. However, this notebook cannot be accessed in the cloud. Users need to either
417 install QIIME2 and GNPS packages on their computer or use Docker. This might be difficult for
418 some, but it is a good option for those familiar with QIIME2⁶⁵. In all cases, users should always
419 consider the size of their data, their computer's power, and their own skill level while using the
420 protocol.

421 **Alternative Open-Source Data Analysis Workflows and Protocols**

422 There have been many efforts in the community to provide and teach statistics solution for non-
423 targeted metabolomics data analysis, and multiple, scripts, web apps and software tools are
424 available for data clean-up, statistical analysis and visualization. While we believe that such a
425 streamlined solution for FBMN results, as described in our protocol, has not yet been provided,
426 we would like to point out the many other tools, workflows and applications that are available.
427

428 **Table 2: Overview of alternative statistics tools and scripting solutions for statistical**
 429 **analysis of non-targeted metabolomics data.**

Tool Name	Tool Type	Availability	Raw Data Processing	Blank Removal	Matrix Transformations	Uni-Variate Statistics	Multi-variate Statistics	Export for Downstream Tools	Customizable	URL	Reference
GUI											
MetaboAnalyst	Web App (GUI)	Open Source	Y	Y	Y	Y	Y	N	N	www.metaboanalyst.ca/	23,24
Workflows											
Galaxy-M	Workflow	Open Source	Y	Y	Y	Y	Y	N	N	github.com/Viant-Metabolomics/Galaxy-M	66
Workflow4Metabolomics	Workflow	Open Source	Y	Y	Y	Y	Y	N	N	github.com/workflow4metabolomics	67
UmetaFlow	Workflow	Open Source	Y	Y	Y	Y	Y	N	N	github.com/biosustain/snake_kemake_UmetaFlow	68
Chemometrics Tutorials	Workflow / Tutorial	Open Source	N	N	Y	Y	Y	N	Y	github.com/Gscoreira89/c_hemometrics-tutorials	
QIIME2 metabolomics plugin	Language	Open Source	N	N	N	Y	Y	N	N	library.QIIME2.org/plugins/q2-metabolomics/10/	65
R Libraries											
mixOmics	Package	Open Source	N	N	Y	Y	Y	Y	Y	mixomics.qfab.org	69
MetaboanalystR	Package	Open Source	Y	Y	Y	Y	Y	Y	Y	www.metaboanalyst.ca/docs/RTutorial.xhtml	70
omu	Package	Open Source	N	N	Y	Y	Y	Y	Y	cran.r-project.org/web/packages/omu/vignettes/Omu_vignette.html	71
metabolomicsR	Package	Open Source	N	N	Y	Y	Y	Y	Y	cran.r-project.org/web/packages/metabolomicsR/index.html	72
MAIT	Package	Open Source	N	N	Y	Y	Y	Y	Y	www.bioconductor.org/packages/release/bioc/html/MAIT.html	73
ropls	Package	Open Source	N	N	Y	Y	Y	Y	Y	bioconductor.org/packages/release/bioc/html/ropls.html	
MSStats	Package	Open Source	N	Y	Y	Y	Y	Y	Y	github.com/Vitek-Lab/MSstats	74
Python Libraries											
TidyMS	Package	Open Source	Y	Y	Y	N	N	Y	Y	github.com/griquelme/tidyms	75

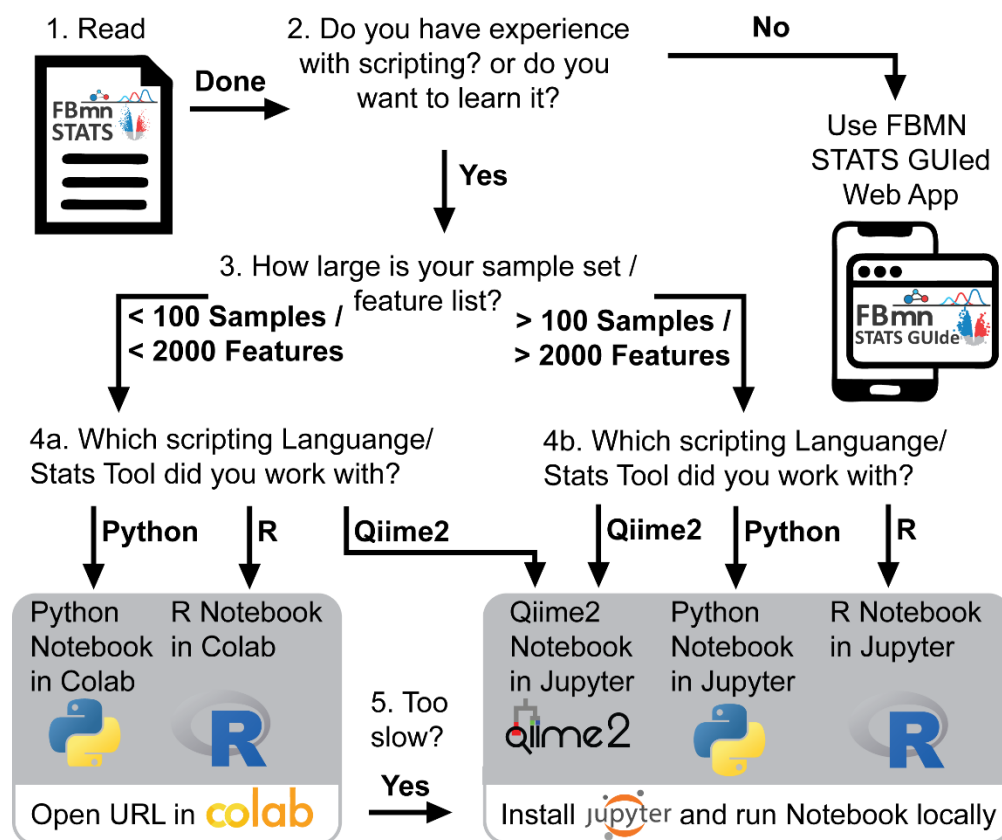
430
 431 We summarized those that, in our opinion, are the most commonly used software tools in **Table**
 432 **2**. This table provides an overview of their functions, purpose, tool type, and when applicable,
 433 references to related protocols and guidelines. We also indicated where in the data processing
 434 pipeline these tools have application by indicating yes (“Y”) or no (“N”) in columns related to raw
 435 data processing (generation of feature quantification table, see **section ‘Feature Detection’**),
 436 data clean-up steps (involving quality control, missing value imputation, normalization, scaling,
 437 and transformation), and multivariate and univariate analyses.

438 We do note that this table is by no means exhaustive. All of these options are workflow dependent
439 and vary based on factors such as the structure of the acquired feature quantification table and
440 the chosen data analysis techniques⁷⁶, and typically require specific file and table formats.

441 Expertise Needed to Implement the Protocol

442 We aimed to make this protocol accessible to a broad range of researchers, from absolute
443 beginners to experts. As we provide different options of executing the code (Web App, Colab and
444 Jupyter notebooks), the protocol should be useful for users both new to metabolomics data
445 analysis, who want to perform a fixed set of processing and statistical analysis, as well as users
446 that require customizable options and need to analyze large datasets. To guide readers through
447 the different options and help to choose which option is most suitable, we generated a decision
448 tree displayed in **Figure 3**. At a minimum, we recommend to have some general background in
449 statistics and a basic understanding of LC-MS/MS data structure, as well as knowledge about the
450 biological system and the experimental design of the dataset which should be analyzed.

451



452

453

454

Figure 3: Decision tree to guide choosing which notebook/app to use.

455 2. Materials

456 Software Used

- 457 ● Google Colab (Optional, cloud)
- 458 ● Local installation of Jupyter Notebook (Optional)
- 459 ● Streamlit Web App (Optional)
- 460 ● QIIME2 (Optional)

461 **Note:** The pipeline can be accessed through Google Colab (which requires no local software
462 installation), a Web App with Graphical User Interface (GUI), or through a local installation of
463 Jupyter Notebook, which may be preferred for larger datasets. A decision tree for which method
464 to choose is provided in **Figure 3**. As a default for beginners, we recommend using the Colab
465 Notebook with R code. In addition to the R code, Python and QIIME2 versions are also available
466 in our GitHub Repository ([https://github.com/Functional-Metabolomics-Lab/FBMN-
467 STATS/tree/main/data](https://github.com/Functional-Metabolomics-Lab/FBMN-STATS/tree/main/data)). Additional information regarding files other than R code available are
468 provided in our supplementary information. A Web App version of the protocol ([https://fbmn-
469 statsguide.gnps2.org/](https://fbmn-statsguide.gnps2.org/)) is available for those who prefer a more visual interface.

470 Required Files:

- 471 ● Feature quantification table (.csv)
- 472 ● Metadata (.txt)

473
474 The feature quantification table (.csv), a characteristic product of LC-MS/MS metabolomics
475 studies, encompasses all mass spectral features (integrated peak areas) and their relative
476 intensities across diverse samples. As mentioned earlier, we used MZmine 3 to obtain the feature
477 quantification table in .csv file format.

478 The metadata is a .txt file that can be created in a word editor or spreadsheet programs such as
479 excel or google sheets. The metadata table needs to be created by the user, providing additional
480 context for the measured samples, such as sample type, species, tissue type, and collection time
481 point. For the datasets to be fully considered for public meta-analysis, we suggest using a
482 standardized metadata format with controlled vocabulary. We recommend adhering to
483 standardized metadata practices and protocols, and for guidance users can refer to the ReDU
484 metadata template (<https://mwang87.github.io/ReDU-MS2-Documentation/HowtoContribute/>).

485 The metadata format in this protocol should be compatible with GNPS workflows ([https://ccms-
486 ucsd.github.io/GNPSDocumentation/metadata/](https://ccms-ucsd.github.io/GNPSDocumentation/metadata/)). The first column in the metadata, labeled
487 'filename', should match the exact filenames as reported in the feature quantification table.
488 Following this, one should include additional columns to the metadata that begin with
489 'ATTRIBUTE_' (e.g., ATTRIBUTE_groups, ATTRIBUTE_timepoint).

490 In our example metadata, columns like ATTRIBUTE_Replicate, ATTRIBUTE_Sample_Type,
491 ATTRIBUTE_Batch, ATTRIBUTE_Month, and ATTRIBUTE_Year all contain group-based
492 information. This type of grouping will assist in selecting different categories for statistical analysis
493 throughout this guide. You can also include columns with continuous numerical data, such as

494 ATTRIBUTE_Injection_order or ATTRIBUTE_timepoint. To ensure statistical power, it is essential
495 to use replicates (we suggest at least three) for each sample type within the experimental design.
496 See **Table 3** for an illustration of the metadata structure.

497

498 **Table 3: Sample metadata layout.**

499 The first column, 'filename', lists the filenames along with their specific extensions (preferably
500 'mzML' or the older 'mzXML'), exactly matching the column names in the feature quantification
501 table. Two example "ATTRIBUTE_" columns are also included: "ATTRIBUTE_groups", which
502 showcases sample categorical data (i.e., different sample types such as Control, Sample, and
503 Blanks), and "ATTRIBUTE_timepoint", which is an example for numerical data.

filename	ATTRIBUTE_groups	ATTRIBUTE_timepoint
control_rep1.mzML	Control	1
⋮	⋮	⋮
Sample_type1_rep1.mzML	Sample_type1	4
Sample_type1_rep2.mzML	Sample_type1	4
⋮	⋮	⋮
blank.mzML	Blank	NA

504 **Additional Input Files**

505 Besides the feature quantification table and metadata, the R and Python notebooks can also
506 integrate molecular annotation files (either in .txt or .tsv format). These include SIRIUS,
507 CANOPUS, and GNPS annotations, which enrich our understanding of each feature during
508 analysis. While the inclusion of SIRIUS and CANOPUS files is optional, they can provide valuable
509 insights.

510 GNPS annotations can be obtained from the Feature-Based Molecular Networking (FBMN)
511 analysis. The process requires MS/MS fragmentation patterns in the ".mgf" format, a feature
512 quantification table, both obtained with e.g., MZmine 3 (see **section 'Feature Detection'**), and a
513 metadata file. The .mgf file carries spectral information about specific MS/MS scans designated
514 for each feature and feature IDs match with feature IDs in the feature quantification table. All of
515 these files need to be uploaded to the GNPS platform.

516 The metabolite annotation requires a user-defined mass tolerance. Subsequently, MS/MS
517 patterns are matched against the GNPS database using a modified cosine similarity⁷⁷, resulting
518 in a molecular network that allows for the identification of compound names for all library hit
519 features. The output of the FBMN job associated with the example data of this protocol is publicly
520 available

521 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=b661d12ba88745639664988329c1363e> and

522 can be downloaded using the ‘Download Cytoscape Data’ option. The FBMN job’s .graphml file,
523 found under the folder “gnps_molecular_network_graphml”, can be used to visualize the
524 molecular network in Cytoscape software. The respective annotated files are located in the
525 “DB_result” and “DB_analog_result” sub-folders (assuming an analog search is performed),
526 with the former offering level 2 and the latter providing level 3 (molecular class) annotations. The
527 analog search identifies structurally related molecules within the molecular network by applying a
528 score threshold, such as a minimum cosine score that MS/MS spectra must achieve to be
529 considered an annotation during spectral matching with MS/MS spectral libraries. An upper limit
530 can be established for the mass shift between the library and potential analogs (e.g., 100 Da),
531 thus expanding the scope of annotation.
532 SIRIUS⁵⁶ can predict molecular formulas, as well as structures through structure database
533 matching using CSI:FingerID^{35,78}. Furthermore, the integrated CANOPUS⁶⁴ module provides
534 ClassyFire based chemical class predictions. As for GNPS, the required input is a .mgf file
535 associated with the MZmine feature quantification table with matching feature IDs across both
536 files. However the .mgf file exported for SIRIUS through MZmine 3 differs from the .mgf exported
537 for GNPS in that it contains isotopic MS1 patterns for accurate molecular formula prediction.
538 All example input files to follow this protocol can be retrieved from the Functional Metabolomics
539 Lab GitHub Repository (<https://github.com/Functional-Metabolomics-Lab/FBMN-STATS/tree/main/data>). Furthermore, users have the convenience of directly uploading all input
540 files by simply entering the task ID from their FBMN job on GNPS.
541

542 Example Dataset

543 The example dataset is part of a previously published study¹⁴, aimed to elucidate the effects of
544 urbanization on organic matter chemotypes in coastal environments after a major rainfall event.
545 Seawater samples were collected from 30 locations over seven areas along the San Diego,
546 California coastline: Torrey Pines, SIO La Jolla Shores (Scripps Institution of Oceanography at
547 La Jolla Shores), La Jolla Cove, La Jolla Reefs, Pacific Beach, Mission Beach, Mission Bay,
548 capturing both pre- (Dec 2017) and post-rainfall (Jan 2018) conditions. In our analysis, we
549 included supplementary data from October 2018, collected from the same sites (no-rain period),
550 for our pipeline evaluation. The dataset consisted of 180 samples from the three sample times
551 (Dec 2017, Jan 2018, Oct 2018) and 2 PPL process blanks at each of the sample times. The
552 datasets can be found in the MassIVE repository: MSV000082312 and MSV000085786
553 <https://massive.ucsd.edu/ProteoSAFe/dataset.jsp?task=8a8139d9248b43e0b0fda17495387756>
554 <https://massive.ucsd.edu/ProteoSAFe/dataset.jsp?task=c8411b76f30a4f4ca5d3e42ec13998dc>

555 **Note:** Seawater samples collected during October 2018 were not available in the original article
556 yet. The .mzML files were preprocessed using MZmine 3 (<https://mzmine.github.io/>) and the
557 feature-based molecular networking workflow in GNPS
558 (<https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=b661d12ba88745639664988329c1363e>).

559 3. Procedure

560 This protocol primarily focuses on the R workflow, given its broad adoption in metabolomics data
561 analysis and the extensive libraries it offers for this purpose. However, recognizing the diversity
562 of our audience and the growing popularity of other platforms, we've also developed workflows in
563 Python and QIIME2 as well as a web application. Please refer to the SI document for the
564 Python/QIIME2 notebook or web app workflow. In the following sections, instead of step-by-step
565 instructions, we highlight the key concepts to avoid repetition from the notebook. Code blocks are
566 included to illustrate the main algorithms and functions.

567 ▲**CRITICAL:** We recommend initially executing the notebook using the provided example
568 dataset. Once familiar, proceed with your own data. This approach ensures a smooth transition
569 from learning to applying the workflow.

570 **General Instructions for Navigating the Notebook:**

571 ● **Text in Red:** These sections indicate critical information or cells that require user input
572 within the notebook. They serve as instructions for adapting the notebook to different
573 datasets without the need to modify the code. Further details are provided within the
574 notebook.

575 ● **User Prompt Guidance:** When you encounter code cells with red highlights, simply
576 execute them without changing their contents. For instance, you may come across lines
577 such as

```
578 Directory <- normalizePath(readline("Enter the folder path in the pop-up  
579 box: "), "/", mustWork=FALSE)
```

580 To provide input, a pop-up box will appear in the output section. Make sure to enter your
581 answers in the pop-up box instead of entering directly within the code. After entering your
582 input, remember to press 'Enter' to proceed to the next step.

583 Using these prompt boxes ensures that user input is seamlessly integrated into the
584 following operations. The position of these prompt boxes might differ depending on your
585 system as they could appear directly below the active cell, at the notebook's top, or even
586 towards the upper section of your screen.

587 ● **Text in Green:** This indicates that the following cell in the notebook contains function
588 definitions and will not display any visible outputs. Even though the underlying code in
589 these cells may seem complex, its purpose is to make repetitive tasks more efficient.
590 Readers who come across these green-highlighted cells do not need to understand the
591 complexities of the code.

592 ● **Using the '#' Operator:** Lines in the code cells that start with '#' are comments explaining
593 the code's function or purpose. These comments are "commented out" and will not be
594 executed. To run a commented-out code, remove the '#' symbol and run the cell again.
595

596 **3.1. Preliminary Setup for the Notebook**

597 We recommend beginners to use Google Colab for the R notebook due to its hassle-free setup
598 as it requires no installations, making it accessible for those unfamiliar with the setup process.
599 However, for regular analysis, local execution in Jupyter on a contemporary desktop computer

600 (E.g., Intel i7, 16 core, 64 GB RAM) is typically faster and more efficient. The reported processing
601 times here are based on our example data set on the Colab platform. The durations other than
602 for package installation are estimated from a beginner's viewpoint, reflecting the time typically
603 required for someone new to complete the analysis. To easily install and run Jupyter Notebook in
604 R, consider using Anaconda Navigator, following the instructions provided in the accompanying
605 document ([https://github.com/Functional-Metabolomics-Lab/FBMN-
606 STATS/blob/main/Anaconda_Rkernel_installation_JupyterNotebook_JupyterLab.pdf](https://github.com/Functional-Metabolomics-Lab/FBMN-STATS/blob/main/Anaconda_Rkernel_installation_JupyterNotebook_JupyterLab.pdf)).

607
608 ▲ **CRITICAL**: To ensure proper execution and chronological order, please run the notebook cell-
609 by-cell instead of running multiple cells simultaneously. The numbers assigned to each cell will
610 help you navigate and determine if the cells have been executed correctly and in chronological
611 order.

612 3.1.1. Package Installation ● Timing 15 mins

613 Step 1: Package Installation

- 614 ● The notebook utilizes R version 4.1.3 (2022-03-10)
- 615 ● Begin by installing and loading the necessary R packages using the `p_load()` function
616 from the 'pacman' (v0.5.1) package⁷⁹. This function checks if a package is installed, if not,
617 it installs the package from CRAN or other repositories in the pacman repository list and
618 loads the package. It is a more efficient alternative to using `install.packages()` and
619 `library()` functions individually for each package.
- 620 ● **Required Packages:** The following R packages are essential for this protocol:
 - 621 ○ Data Cleanup: tidyverse⁸⁰ (v2.0.0), IRdisplay⁸¹ (v1.1), KODAMA⁸² (v2.4)
 - 622 ○ Multivariate Statistics: factoextra⁸³ (v1.0.7), vegan⁸⁴ (v2.6-4), ComplexHeatmap⁸⁵
623 (v2.10.0), dendextend⁸⁶ (v1.17.1), NbClust⁸⁷ (v3.0.1), rfPermute⁸⁸ (v2.5.1).
 - 624 ○ Univariate Statistics: FSA⁸⁹ (v0.9.4), matrixStats⁹⁰ (v0.63.0).
 - 625 ○ Visualization and Plotting: ggsci⁹¹ (v3.0.0), cowplot⁹² (v1.1.1), svglite⁹³ (v2.1.1).
- 626 ● Packages are installed just before their respective sections to reduce installation time.
627 However, please note that the packages installed initially in one section can be used for
628 the later sections as well. For example, tidyverse (v2.0.0) can be used throughout the
629 notebook, not just for data cleanup.

631 Step 2: Set Working Directory

632 (User Input Required)

- 633
- 634
- 635 ● Set a folder as the working directory. This is where you access input files and save output
636 files. Make sure to include all necessary input files in this folder.

- 637
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- In Google Colab, click on the three dots in the upper left corner to see the notebook contents. Click on the folder icon and create a new folder by right-clicking in the empty space and selecting 'new folder'.
 - When you run the following cell in the notebook to set a working directory, a pop-up box will display as shown in **Figure 4**. Insert the path of the folder containing your input files and press 'Enter'.
 - **(For Local Environment)** If you're running the notebook in your local environment, you can directly specify the local path of your folder to set it as your working directory. For example, if your folder is located at D:\User\Project\Test_Data, simply input this path when prompted and press 'Enter'.

Directory <- normalizePath(readline("Enter the folder path in the output box: "),"/",mustWork=FALSE)
setwd(Directory)

... Enter the folder path in the output box:

647

648 **Figure 4: Screenshot of the code cell from R Google Colab Notebook to set the working directory.**

649 The image displays the code cell targeting the '/content/test_data' directory. This user-created directory

650 holds the input files for the data analysis. Note the stop symbol with the surrounding loading circle, indicating

651 the cell awaits user input. To proceed to the next cell, provide the input (e.g., /content/test_data) and press

652 enter.

653 3.1.2. Data Import

654 This section guides users through the process of importing necessary data files for the Notebook.

655 Various methods are outlined, catering to different data sourcing preferences.

657 **Step 3: Uploading Files to Google Colab**

658 *(User Input Required)*

659

660 Right-click on the folder you created in the Google Colab workspace and select 'upload' to transfer

661 the required files from your local machine to the cloud session. If you do not want to use files from

662 your local machine, you can skip this step and proceed directly to step 4.2 ('Loading Files from

663 URL') or 4.3 ('Loading Files from GNPS').

664 **Step 4: Select a Data Loading Method (Choose One Option from Steps 4.1 to 4.3)**

665 The user can choose from the following options provided in steps 4.1 to 4.3 to import their data.

667 **Step 4.1: Loading Files from the folder**

668 *(User Input Required)*

669

670

a)

INDEX	FILENAMES
<int>	<chr>
1	20221102_SD_BeachSurvey_batchFile.xml
2	20221125_Metadata_SD_Beaches_with_injection_order.txt
3	GNPS_analog_result_FBMN.tsv
4	GNPS_result_FBMN.tsv
5	SD_BeachSurvey_GapFilled_quant.csv

b)

```
input_str <- readline('Enter the index number of the feature table and metadata separated by commas: ')
Enter the index number of the feature table and metadata separated by commas:
5,2
```

671
672 **Figure 5: Screenshots illustrating loading input files from a folder:** a) Table showing all the files in the
673 working folder, where the first column, labeled “INDEX”, denotes the serial or **index number** of the files.
674 b) Shows the user input interface. Upon executing the code cell, the user is prompted to enter the index
675 numbers for the feature table and metadata. In this example, “5” and “2” are entered, referencing the files
676 indicated in (a).

677
678 In this step, you will begin by viewing a table displaying a list of files in your working folder (e.g.,
679 uploaded in the previous step) as shown in **Figure 5a**. Each file will have an index number
680 associated with it. Your task will be to import three tables by specifying the index number
681 associated to each: the feature quantification table (**ft**), the metadata table (**md**), and optionally,
682 annotation tables (**an**). For an example, see **Figure 5b**. To guide you through this process, there
683 are three code blocks that require user input.

- 684
- 685 1. **Feature Quantification Table and Metadata Import:** The first code block will prompt you
686 to enter the index numbers associated with the feature quantification table and metadata,
687 separated by commas. Simply input the corresponding index number assigned to each of
688 these files.
 - 689 2. **Annotation Tables Import:** The second code block will request the index numbers
690 associated with the annotation tables. Specifically, you will be asked to enter the index
691 numbers of the GNPS library annotation file and the analog annotation files. If you have
692 not performed an analog search for FBMN, only provide the index number of the GNPS
693 library annotation file.
 - 694 3. **SIRIUS Annotation File Import (Optional):** The third code block requests you to input
695 the index number of a SIRIUS annotation. This file will be used to merge all annotations
696 (e.g., GNPS library, analog hits, SIRIUS) into a single master table for easier data
697 exploration later on. It is worth noting that this protocol does not specifically focus on
698 SIRIUS annotations for analysis. The inclusion of SIRIUS annotations is solely for the
699 convenience of consolidating all annotations in one place for the user.

700
701 By following these prompts, one can successfully load the essential tables required for the
702 subsequent analysis. Make sure to carefully input the correct index numbers.

703 **Step 4.2: Loading Files from URL**

704 *(User Input Required)*

- 705
- 706 ● We also provide an example of retrieving data from a URL (for example, the feature
707 quantification table can be obtained from [https://raw.githubusercontent.com/Functional-
708 Metabolomics-Lab/FBMN-STATS/main/data/SD_BeachSurvey_GapFilled_quant.csv](https://raw.githubusercontent.com/Functional-Metabolomics-Lab/FBMN-STATS/main/data/SD_BeachSurvey_GapFilled_quant.csv))
 - 709 ● We access the feature quantification table, metadata, and analog result files directly from
710 our Functional Metabolomics GitHub page.
 - 711 ● If you are using your own dataset (or the test dataset) in Google Colab, you can get the
712 file URL by uploading the input files to the Colab workspace, right-clicking on the file,
713 selecting “Copy path”, and then replacing the URL in the relevant cell.
- 714

715 **Step 4.3: Loading Files from GNPS**

716 *(User Input Required)*

- 717
- 718 ● In this step, you can load files directly from the repositories MassIVE or GNPS. If you have
719 performed FBMN on your feature quantification table, you can access the required files
720 by providing the task ID.
 - 721 ● To locate the task ID of your FBMN job within your GNPS account, navigate to the ‘Jobs’
722 section. Here, the ‘unique ID’ for each job is listed in the ‘Description’ column.
 - 723 ● When you run the relevant cell in the notebook, you will be prompted to enter the task ID
724 within the notebook. Given the task ID, the notebook will retrieve the necessary files for
725 further analysis.
- 726

727 ▲ **CRITICAL:** Make sure your metadata has the necessary attribute columns to describe the data
728 (at least one, e.g., ATTRIBUTE_SampleType). If your FBMN metadata is insufficient, you might
729 need to load additional metadata from a local folder for downstream statistical analysis, adding
730 an additional step in the workflow.

731

732 **Step 5: Exploring the Imported Files**

- 733
- 734 Use the `head()` and `dim()` functions to get an initial view of your imported data files.
- 735 ● The `head(ft)` function displays the first six rows of the feature table by default, giving
736 you a quick look at your data’s structure.
 - 737 ● The `dim(ft)` function reveals the dimensions of your feature quantification table, i.e., the
738 number of rows and columns.
- 739

740 ▲ **CRITICAL:** If you encounter an error while executing certain code cells, it is good practice to
741 verify the correctness of your data tables using the `head()` and `dim()` functions.

742

743 We also provide a special summary function `InsideLevels(md)` to explore the metadata, which
744 returns a summary table with columns for INDEX, ATTRIBUTES, LEVELS, COUNT, and
745 ATTRIBUTE_CLASS.

- 746
747
748
749
750
751
752
753
754
1. **INDEX:** Row number in the summary table
 2. **ATTRIBUTES:** Column name of the attribute, e.g., ATTRIBUTE_Sample_Type
 3. **LEVELS:** Unique groups within the attribute column, e.g., Blanks, Sample
 4. **COUNT:** Number of files for each category, e.g., 6, 180 indicating 6 files for “Blank” sample type and 180 for “Sample” sample type.
 5. **ATTRIBUTE_CLASS:** Data type of the attribute. Useful for spotting cases where a numeric attribute like ATTRIBUTE_minutes is classified as ‘character’.

755 3.1.3. Merging Annotations with Feature Quantification Table

756 This section involves integrating various annotations, such as SIRIUS and GNPS annotations,
757 into our feature quantification table. This process is vital as it helps identify the metabolites
758 corresponding to the features in our feature quantification table, aiding in the interpretation of our
759 metabolomics data.

760 Step 6: Identifying Appropriate Columns for Merging

761 Depending on the type of annotation to be merged, the feature quantification table’s unique ‘row
762 ID’ column is matched with the corresponding column in the annotation file:

- 763
764
- 765 ● **GNPS Annotations:** The ‘row ID’ is matched with the ‘#Scan#’ column in the GNPS
766 annotation file. The ‘Compound_Name’ column contains the annotation information.
 - 767 ● **GNPS Analog Annotations:** Similar to GNPS annotations, the ‘row ID’ is matched with
768 the ‘#Scan#’ column in the GNPS analog annotation file. The ‘Compound_Name’ column
769 contains the annotation information.
 - 770 ● **SIRIUS Summary Files:** The ‘row ID’ is matched with the ‘id’ column in the SIRIUS
771 summary file. A typical feature ID in the ‘id’ column might look like this:
772 “3_ProjectName_Mzmine 3_SIRIUS_1_16”, where the last string (16) represents the row
773 ID.

774 Step 7: Ensuring Data Compatibility

775 Before merging, we ensure that the classes (or data type) of the columns meant to be merged
776 are the same. Then, we can combine feature and annotation data based on the appropriate
777 matching columns. Any mismatch, such as one column being of character type while the other
778 one is numeric, can cause merge errors, even if the values within the columns are identical.

780 Step 8: Merging Annotations

- 781
- 782 ● Rename the column names of analog annotation dataframe ‘an_analog’ with an
783 ‘Analog_’ prefix and merge the modified ‘an_analog’ dataframe with ‘an_gnps’ based
784 on #Scan#.
 - 785 ● For each unique ‘#Scan#’, consolidate multiple compound names into a single row. If both
786 the GNPS compound names (actual and library hits) for a particular ‘#Scan#’ are identical,

787 keep one; otherwise, combine them using a “,” separator. The output is
788 ``an_final_single``.

- 789 • Merge ``an_final_single`` with the feature quantification table (``ft``) using ``#Scan#``
790 and ``row ID`` as matching columns respectively. Keep all rows from the feature
791 quantification table.

792

793 **Additional steps:**

794

795 **Incorporating Additional Annotations (optional)**

- 796 • If SIRIUS annotations are available, follow these additional steps: Extract ``row ID`` from the
797 ``id`` column in the SIRIUS dataframe, rename the columns with a ``SIRIUS_`` prefix, and
798 merge the modified SIRIUS dataframe with ``an_final`` data frame based on ``row ID``.
- 799 • For simplicity, we have shown here how to merge SIRIUS summary files. This process
800 can be similarly adapted for CANOPUS summary files.

801

802 **Exporting the Merged Annotations**

- 803 • Finally, write the merged annotation table to a CSV file for further data exploration and
804 downstream analyses.

805

806 **3.1.4. Ensuring Metadata and Feature Quantification Table** 807 **Compatibility for Downstream Analysis**

808 This section streamlines the metadata and feature quantification tables, ensuring they align
809 perfectly for subsequent steps in the protocol and remove discrepancies between them. By
810 following the outlined steps, we achieve harmonized data structures. A final verification confirms
811 that all files in the feature table are mirrored in the metadata, and vice versa. Upon successful
812 validation, the tables are set for the next section of analysis. If there is a mismatch, often due to
813 naming inconsistencies or missing files, the user needs to rectify these issues before moving
814 forward. As a user, you are not required to modify any of the code within this section. Simply
815 execute each cell in turn.

816

817 **Step 9: Creating Backup Files**

818 The feature quantification table (``ft``) and metadata (``md``) files are stored under different names
819 (``new_ft``, ``new_md``) to preserve the original versions.

820

821 **Step 10: Cleaning up the Feature Quantification Table**

- 822 • Clean the feature quantification table by removing ``peak area`` extensions from the column
823 names, a default format included in MZmine-derived feature quantification tables.
- 824 • Check and remove any columns containing only NA values present in the feature and
825 metadata tables.
- 826 • Check and remove any rows and columns containing only empty strings in the metadata
827 table

828

829 **Step 11: Updating the Row Names of the Feature Quantification Table**

- 830 ● In this step, we reformat the row names to consolidate essential information about each
831 feature. By doing this, we can retain only the numeric data in the feature quantification
832 table and remove all other columns.
- 833 ● The row names are constructed by concatenating the Feature ID, *m/z*, RT, and GNPS
834 annotations into a single string, in the following format:
835 ``XFeatureID_m/z_RT_GNPS_annotations``.
- 836 ● An example row name is ``X92649_226.951_14.813_NA;TRYPTOPHAN``. Here, “NA”
837 indicates that there was no direct library hit for this feature. However, the analog
838 annotation suggests it could be tryptophan.
- 839 ● In the R environment, a dataframe’s row names must be characters or strings. Thus, we
840 add the ‘X’ character prefix to the numeric Feature ID to ensure compatibility.

841

842 **Step 12: Selecting Relevant Columns**

843 *(User input - Optional)*

- 844 ● Retain only ‘.mzML’ (or ‘.mzXML’) file-relevant columns and remove extraneous
845 information, such as additional columns added due to IIMN. Here, the features are
846 represented as rows in the feature quantification table.
- 847 ● Only when the file extensions ‘.mzML’ or ‘.mzXML’ are not available, the user is prompted
848 to enter their respective file extension.
- 849 ● This step ensures that the feature quantification table contains only the intensity values of
850 the features, which is crucial for subsequent calculations. The modified row names provide
851 basic feature information, and for a more detailed understanding, you can refer back to
852 the original feature quantification table.

853

854 **Step 13: Verifying File Consistency**

855 The metadata and feature quantification tables are arranged in the same order of ‘.mzML’ (or
856 ‘.mzXML’) file names. We then verify consistency between the feature and metadata tables by
857 using the ``identical(new_md$filename, colnames(new_ft))`` command.

- 858 ● If the result is TRUE, proceed to data cleanup.
- 859 ● If FALSE, there might be missing files or discrepancies in file naming. Check the
860 corresponding column names in the feature quantification table for potential errors like
861 spelling mistakes or case-sensitive issues, and re-upload the correct files. Re-run all the
862 above steps once corrected.

863

864 **3.2. Data Cleaning: ● Timing 20-30 mins**

865 Following the LC-MS/MS data pre-processing with MZmine, we perform the post-processing of
866 the data (also known as data pretreatment or data clean-up) as the first crucial step in our
867 workflow. While the ‘preliminary setup for the notebook’ section prepares the feature and
868 metadata tables for analysis, actual modifications to the data commence from this section.

869

870

871

872 **Step 14: Transposing the Feature Quantification Table**

- 873 • As a first step, we transpose the feature quantification table. The result is a table (`ft_t`)
874 where the row names represent the sample names, and the column names consist of
875 concatenated feature information.
- 876 • Then, we merge this transposed feature quantification table (`ft_t`) with the metadata
877 (`new_md`), using the sample names as the common link. This merged table, referred to
878 as `ft_merged`, consolidates all necessary information in a single structure.
- 879 • The `ft_merged` table can also be exported to a CSV file for future use, such as batch
880 correction or other specialized analyses.

881

882 **3.2.1. Batch Correction (Optional)**

883 Batch effects are systematic differences in sample measurements when samples are run as
884 multiple batches or groups. In most cases, when the sample sizes exceed the measurement
885 assay, it is often necessary to measure the samples in multiple batches. This might lead to varying
886 mass spectra among the samples within a batch and among different batches⁹⁴. Several factors
887 could contribute to these batch effects such as variability in instrument conditions, RT shifts, and
888 gradual contamination of LC columns when measuring multiple samples over a long period.
889 These are often unavoidable issues, hence it is common to treat these effects post-sample
890 measurement⁹⁵.

891 To correct these unwanted variations, we first need to identify their presence, remove or adjust
892 the variations for further statistical analysis and assess the performance of our method⁹⁶. The
893 most common method for visualizing or identifying the presence of batch effects is through a
894 simple Principal Component Analysis (PCA), guided PCA⁹⁷ or Principal Coordinates Analysis
895 (PCoA). In the PCA/PCoA scores plot, it is generally expected that all the QCs cluster together
896 indicating little analytical variation in the data. When the inter-batch variation gets higher, the inter-
897 QC distances in the PCA/PCoA plot will also increase⁹⁸. To visually assess this using the
898 notebook, follow these steps:

899 1. Execute Step 25 to install necessary packages for multivariate analysis.

900 2. Run Step 32 and Step 33 to visualize the PCoA using the custom-made `plotPCoA()`
901 function. Detailed usage instructions are provided in the respective steps.

902 Assuming your sample type information (description of which samples are pooled QCs, blanks,
903 samples etc.) is located in the `ATTRIBUTE_Sample_Type` column of the metadata, the function
904 can be invoked as follows:

```
plotPCoA(  
  ft = ft_t,  
  md = new_md,  
  distmetric = "euclidean",  
  category_permanova = "ATTRIBUTE_Sample_Type",  
  pcoa_category_type = 'categorical',  
  category_pcoa_colors = "ATTRIBUTE_Sample_Type")
```

905 We have deferred this visualization to a later section, after data cleanup. As we delve deeper into
906 multivariate analyses after data cleanup, this approach avoids redundancy and ensures users
907 can maximize the utility of this protocol.

908 Another method is to use Analysis of variance (ANOVA) by comparing the QC mean of different
909 batches for statistically significant differences⁹⁹. Once the presence of a notable batch effect is
910 confirmed, multiple approaches can be used to correct the effects, including

- 911
- 912 • **Normalization** methods such as Metabodrift¹⁰⁰, ComBat¹⁰¹
- 913 • **Transformation** method: wavelCA¹⁰² (wavelet transformation coupled to ICA)
- 914 • **Regression-based** approaches such as linear least-square (LS) method¹⁰³, QC-based
915 robust LOESS correction⁴⁰, QC-support vector regression¹⁰⁴ (QC-SVR)
- 916 • **ML-based** methods such as random forest-based QC-RFSC correction¹⁰⁵, deep learning
917 model: NormAE (Normalization Autoencoder) algorithm¹⁰⁶, Regularized Adversarial
918 Learning Preserving Similarity¹⁰⁷ (RALPS).

919 Each method has its strengths and limitations. When there are no QCs included in the study,
920 normalization can be used instead to attempt to reduce most of the unwanted variations¹⁰⁸, at the
921 risk of removing true biological variation. For the sake of simplicity and to cater primarily to
922 beginners, this protocol does not elaborate on batch correction. However, for those interested in
923 exploring batch correction in depth, we have prepared a supplementary R notebook available on
924 our GitHub repository (https://github.com/Functional-Metabolomics-Lab/FBMN-STATS/blob/main/R/Additional_Notebooks/Batch_Correction.ipynb). In this notebook, we
925 execute inter-batch correction similar to the method described by Qin Liu *et al*⁹⁴. The procedure
926 involves calculating the mean of each feature across all batches, then calculating the batch-
927 specific feature mean, and subsequently adjusting feature intensities within each batch relative to
928 the batch-specific and overall means. For intra-batch adjustments, the notebook illustrates the
929 QC-based LOESS correction method, with a prerequisite that each batch should start and end
930 with a pooled QC injection.
931

932 3.2.2. Blank Removal

933 To prioritize or identify metabolites from our samples, we need to remove contaminants, i.e.,
934 features found in the blanks, before proceeding with statistical analysis¹⁰⁹. While blank removal
935 can be executed during pre-processing with MZmine 3, which might result in the absence of blank
936 features and samples in both the feature table and metadata, conducting it during post-processing
937 offers more flexibility. If you have performed blank removal during pre-processing, simply skip
938 steps 16-18. Instead, designate the previous variables as the resulting blank-removed table and
939 the metadata for samples, ensuring a seamless continuation of the workflow: `blk_rem <-
940 ft_t` and `md_Samples <- new_md``. For a graphical overview on blank removal, see **Figure**
941 **6**, and for more insights, refer to **Box 2**.

942
943
944

945 **Step 15: Examining Metadata Attributes**

946 Run `InsideLevels(new_md)` to identify unique groups within each metadata attribute. This
947 helps to find the attribute column containing sample type information (e.g., 'Blanks', 'Samples').

948

949 **Step 16: Separating Blank and Sample Files**

950 *(User Input Required)*

951

952 In this step, the data is split into two groups: 'blank' and 'sample' files. It's important to note that
953 'samples' here include all mzML (or mzXML) files except blanks, including control samples, as
954 they might be influenced by blank features.

955

956

- **Identify the Attribute Column:** The user will first be prompted to enter the index number
957 of the attribute containing information about samples and blanks. Here, it is
958 `'ATTRIBUTE_Sample_Type'`.

959

- **Display Unique Groups:** The unique groups within the chosen attribute column will be
960 displayed. For example, in our dataset, it will show Blank and Sample. However, your
961 dataset might include various groups, such as Blank, Samples, Control, etc.

962

- **Select the Blank group:** Next, the user will be prompted to enter the index number
963 corresponding to the blank group. If there are multiple groups representing blanks (e.g.,
964 Blank, PPL_Blank), their index numbers should be entered, separated by commas.

965

- **Select the Sample group:** Similarly, the user will be asked to enter the index number(s)
966 for the sample level. If the dataset includes multiple groups for samples (e.g., Sample,
967 Control), the corresponding index numbers should be entered, separated by commas.

968

- **Subset the Data:** Using the information provided, the metadata (`'new_md'`) will be
969 subsetted into `'md_Blank'` and `'md_Samples'`. The corresponding feature
970 quantification tables will be obtained and named `'Blank'` and `'Samples'`, respectively.

971 **Step 17: Define Cutoff for Blank Feature Removal**

972 *(User Input Required)*

973

974 In this step, the user will need to set a cutoff value within the range of 0 to 1, with a recommended
975 range of 0.1 to 0.3. This value will determine which features are considered to be artifacts of the
976 blank and thus removed from the dataset. The next step will explain how the features exceeding
977 this cutoff are identified and eliminated.

978

979 **Step 18: Perform Blank Removal**

980 Calculate the blank's contribution to each feature and eliminate those exceeding the user-defined
981 cutoff. This is achieved by:

982

- Compute the mean value for each feature within the dataframes (`'Blank'`) and
983 (`'Samples'`). This step calculates two mean values for each feature, one for blanks and

984 one for samples. These averages are stored in a new dataframe called ``Avg_ft`` under
985 the columns ``Avg_blank`` and ``Avg_samples``.

986 • Compute the ratio of the average blank contribution to the average sample contribution
987 for each feature.

988 • Generate a binary mask where entries corresponding to ratios above the user-defined
989 cutoff are marked as 1 (TRUE), and all others are set to 0 (FALSE). This mask helps in
990 identifying which features are significantly present in blanks as compared to samples.

991 • Retain only the features associated with 0s in the binary mask. Features with a ratio
992 exceeding the cutoff (marked as 1) are considered artifacts from the blanks and are thus
993 removed. Conversely, if the feature intensity is significantly higher in samples than in
994 blanks, it is deemed a true feature from the samples and is retained (marked as 0).

995 • The final table, free from blank artifacts, is named ``blk_rem``, and its corresponding
996 metadata is ``md_Samples``.

997 The final output is the ``blk_rem`` table, which excludes background or noise features.
998 Information on the total number of features, the number of background/noise features, and the
999 number of features after noise exclusion is also displayed. Steps 16-18 are displayed in **Box 2**.

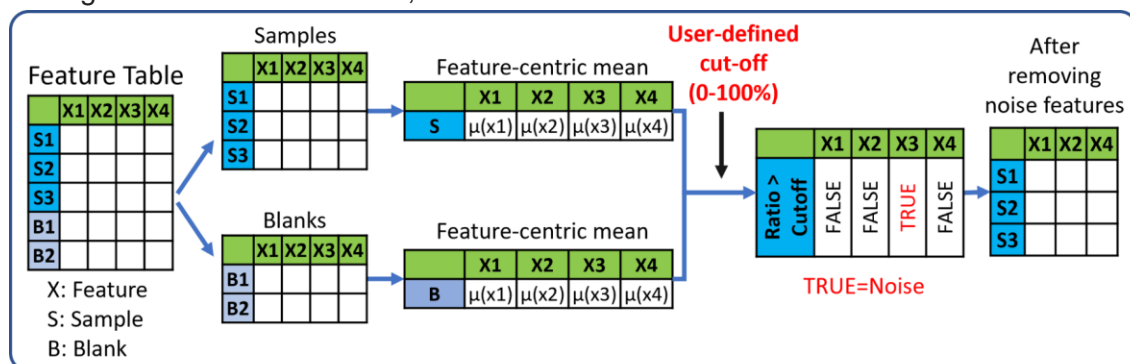
1000
1001 ▲ **CRITICAL:** Lowering the cutoff to 0.1 demands a greater contribution from the sample (90%)
1002 and limits the blank's contribution to 10%. Raising the cutoff leads to fewer background features
1003 being identified and more analyte features being observed. Conversely, lowering the cutoff is
1004 more stringent and removes more features.

1005 1006 **Box 2 - Blank Removal**

Some existing methods to achieve blank removal are: creating a molecular network using the online platform, the global natural product social molecular networking (GNPS), and visualizing the network in Cytoscape to manually remove the blank and media nodes. But this is a tedious process⁷; there is also Lawson *et al.*'s msPurity R package with a function called "SubtractMZ" to perform blank removal¹¹⁰. Data-adaptive filtering methods have also been suggested to remove features from blanks and low abundant features from samples with undetected values¹¹¹.

Another popular feature filtering method is based on the Coefficient of Variance (CV). Also referred to as relative standard deviation (RSD) is a measure of statistical dispersion, calculated as the ratio of the standard deviation to the mean¹¹². When pooled QC samples are integrated throughout a study, CV can be used to assess the stability of each feature. As a general rule of thumb, features exhibiting a CV greater than 30% are typically excluded, though the threshold is more stringent (at 20%) for FDA studies. However, it's essential to approach CV filtering with caution. Schiffman *et al.* have highlighted the potential limitations of this method, pointing out that CV primarily evaluates variability across technical replicates without giving weight to biologically meaningful variability across different subjects¹¹³. Consequently, while CV filtering

might be apt for studies focusing on homogenous samples like plasma or *Escherichia coli* cells, it might not be the best fit for diverse sample sets such as environmental or fecal samples. The dispersion ratio or D-ratio, introduced by Broadhurst *et al.*, offers an alternative to a simple CV cut-off by comparing both technical and biological variance. It is calculated by dividing technical variance by the total variance, which includes both technical and biological variances. Therefore, for any feature, a 0% D-ratio signifies that the variance is entirely biological, whereas a 100% D-ratio denotes complete technical noise, without any biological information. So, when assessing D-ratios for metabolites, it is better to retain the ones with D-ratios closer to zero¹¹⁴.



```
#Getting mean for every feature in blank and Samples in a data frame named 'Avg_ft'

Avg_ft <- data.frame(Avg_blank=colMeans(Blank, na.rm= F))
# set na.rm = F to check if there are NA values. When set as T, NA values are changed to 0

# adding another column 'Avg_samples' for feature means of samples
Avg_ft$`Avg_samples` <- colMeans(Samples, na.rm= F)

#Getting the ratio of blank vs Sample
Avg_ft$`Ratio_blank_Sample` <- (Avg_ft$`Avg_blank`+1)/(Avg_ft$`Avg_samples`+1)

# Creating a bin with 1s when the ratio>Cutoff, else put 0s
Avg_ft$`Bg_bin` <- ifelse(Avg_ft$`Ratio_blank_Sample` > Cutoff, 1, 0 )

#Calculating the number of background features and features present
print(paste("Total no.of features:",nrow(Avg_ft)))
print(paste("No.of Background or noise features:",sum(Avg_ft$`Bg_bin` ==1,na.rm = T)))
print(paste("No.of features after excluding noise:",(ncol(Samples) - sum(Avg_ft$`Bg_bin` ==1,na.rm = T))))

blk_rem <- merge(as.data.frame(t(Samples)), Avg_ft, by=0) %>%
  filter(Bg_bin == 0) %>% #picking only the features
  select(-c(Avg_blank,Avg_samples,Ratio_blank_Sample,Bg_bin)) %>% #removing the last 4 columns
  column_to_rownames(var="Row.names")
blk_rem <- data.frame(t(blk_rem))
```

```
[1] "Total no.of features: 11217"
[1] "No.of Background or noise features: 2125"
[1] "No.of features after excluding noise: 9092"
```

```
head(blk_rem, 2)
dim(blk_rem)
```

```

X10015_282.169_2.763_NA X10035_325.139_2.817_NA X10037_216.123_2.847_NA X10047_338.159_2.845_NA X10058_280.117_2.961_NA
<dbl> <dbl> <dbl> <dbl> <dbl>
SD_01-
2018_1_a.mzXML 50907.97 196008.38 90480.91 446560.7 182757.8
SD_01-
2018_1_b.mzXML 51443.73 99569.05 411595.38 239022.0 274146.0
180 · 9092
```

Figure 6: Blank Removal Process: Featuring a graphical representation of the blank removal followed by screenshots of the corresponding R code executed for the procedure.

1007

1008 3.2.3. Imputation

1009

1010 Many feature extraction software programs, such as MZmine 3, often generate tables with missing
1011 values denoted as “NA”, “NaN” or 0. This means that for several *m/z* and RT traces in a given
1012 sample, there may not be a peak detected and therefore no value is available⁷⁶. However, many
1013 statistical approaches, such as Principal Component Analysis (PCA), require numerical values
1014 for each observation. Hence, these features with missing values need to be removed or imputed.
1015 In this section, we handle the zero values in our blank-removed feature quantification table. Refer
1016 to **Figure 7** for a graphical overview on imputation and for more information, see **Box 3**.

1017

1018 **Step 19: Analyzing the frequency distribution of relative intensities**

1019 We first examine the distribution frequency of the relative intensities in our feature quantification
1020 table by creating a histogram. This reveals any notable gaps in the range of values, such as a
1021 large number of zeros or a lack of values within a particular range. In our example, we observed
1022 many zeros and no values in the range of 0 to 100. The smallest non-zero value in our table was
1023 between 100 and 1000.

1024

1025 **Step 20: Replacing zeros with random values**

1026 We replace all zero values in the dataset with the randomly generated number between 1 and the
1027 smallest non-zero value in our blank-removed table. This process, known as imputation, fills in
1028 the gaps in our data with plausible values, which can improve subsequent analyses.

1029 ▲ **CRITICAL:** Imputation is not advised if one plans to execute a PCoA using the Jaccard distance
1030 since Jaccard transforms data into binary (0 and 1). Without zeros, it results in a table full of ones.

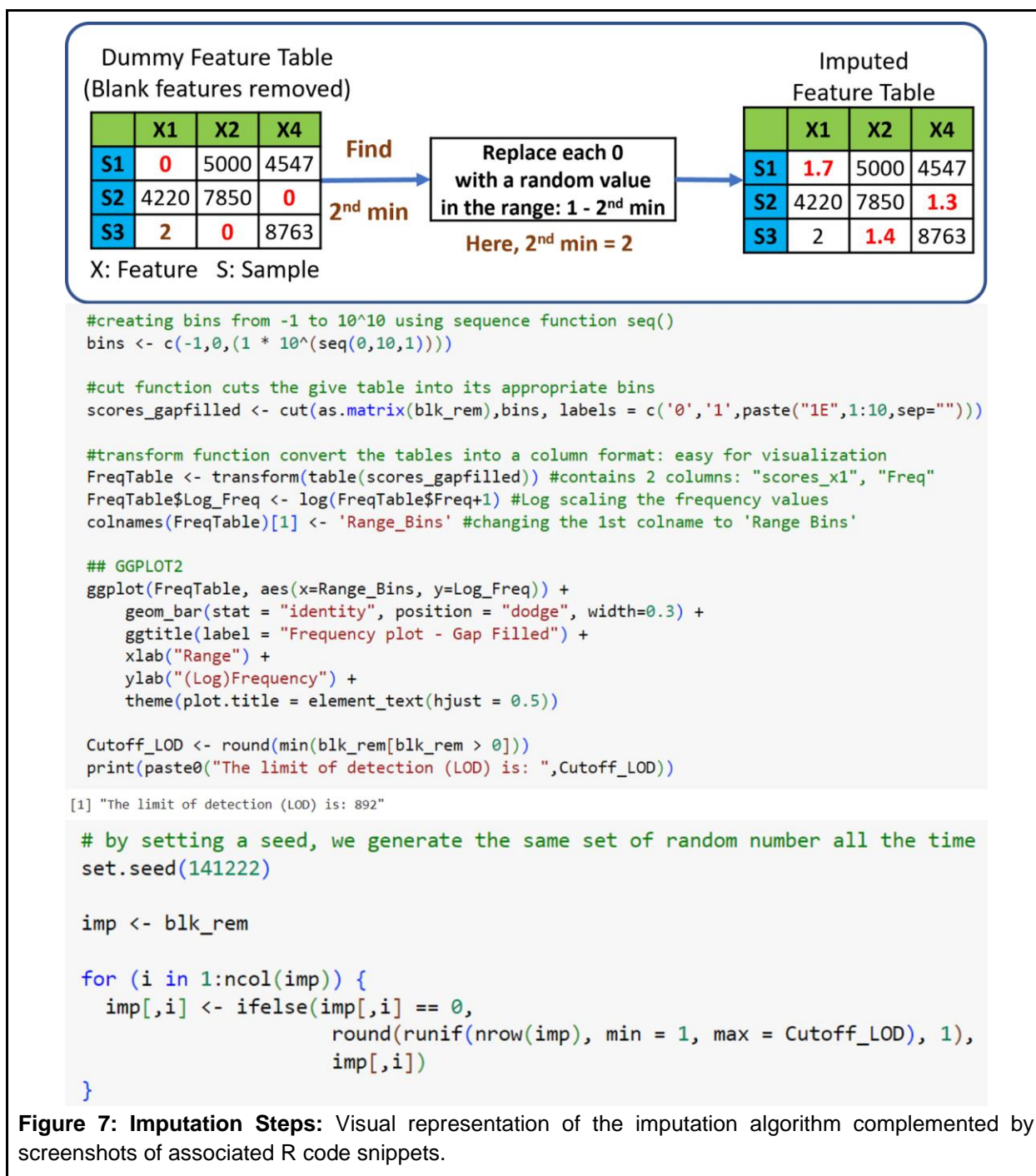
1031

1032 **Box 3 - Imputation strategies**

The appropriate imputation strategy depends on the nature of the missing values:

1. Below the Limit of Detection (LOD): If a value is missing because the corresponding molecule was below the analytical method's LOD, consider replacing missing values with a low value, ensuring it does not artificially lower the variance¹¹⁵. Our imputation method corresponds to this scenario.

2. Sample Processing or Feature Extraction Artifacts: If missing values arise from analysis anomalies (like ion suppression effects or specific retention time shifts) or sample processing artifacts, consider substituting missing values with those similar to values detected in other samples. Here, machine learning methods like k-nearest neighbor (KNN) or random forest (RF) can be useful. KNN fills in multiple missing values by identifying the k nearest data points to a given point¹¹⁶.



1033

1034 3.2.4. Normalization

1035 Sample normalization aims to eliminate systematic bias via adjusting variations across
 1036 samples¹¹⁷. In our pipeline, we show two normalization methods: Total Ion current (TIC)
 1037 normalization and Probabilistic Quotient Normalization (PQN), implemented using the KODAMA
 1038 library in our R Notebook. Therefore, we begin this section by installing the KODAMA package.

1039 We recommend that users run both normalization methods and scaling methods (steps 21 to 23),
1040 but they can choose either method for further analysis in step 24. Additional information about
1041 normalization, including various methods, and guidelines for selecting the most suitable method
1042 for a given dataset, is provided in the accompanying **Box 4**. For a graphical view on the provided
1043 normalization methods, see **Figure 8**.

1044 **Step 21: Total Ion Current (TIC) Normalization**

1045
1046 In TIC normalization, also known as total sum normalization, every feature within a sample is
1047 normalized relative to the area of the TIC chromatogram¹¹⁸. This involves dividing each feature
1048 by the sum of peak areas of all features within a sample. The normalization function from the
1049 KODAMA package performs row-wise sum operations, and we have the samples arranged in
1050 rows and features in columns.

1051 1052 **Step 22: Probabilistic Quotient Normalization (PQN)**

1053 PQN is another method performed on the imputed table, resulting in a PQN-normalized table with
1054 features in columns and samples in rows.

1055
1056 PQN is based on the comparison of a ‘test’ spectrum (the individual sample to be normalized)
1057 with a ‘reference’ or ‘control’ spectrum. The steps involved in PQN are as follows¹¹⁹:

- 1058 ● **Normalization of Test Spectrum:** The test spectrum is first normalized, typically using a
1059 sum normalization technique like TIC.
- 1060 ● **Selection of Control Spectrum:** The control spectrum acts as a standard for comparison.
1061 It could be a pre-determined standard obtained from a database or calculated as the mean
1062 or median spectrum from all samples or quality control (QC) samples.
- 1063 ● **Calculation of Quotients:** For each sample, quotients are calculated between the
1064 features in the test spectrum and the corresponding features in the control spectrum. This
1065 step results in a median quotient spectrum for each sample.
- 1066 ● **Normalization by Median Quotient Spectrum:** Each test spectrum is then normalized
1067 by dividing it by its corresponding median quotient spectrum. This process scales the test
1068 spectrum values relative to the control spectrum, ensuring an equal basis for comparison
1069 across all samples.

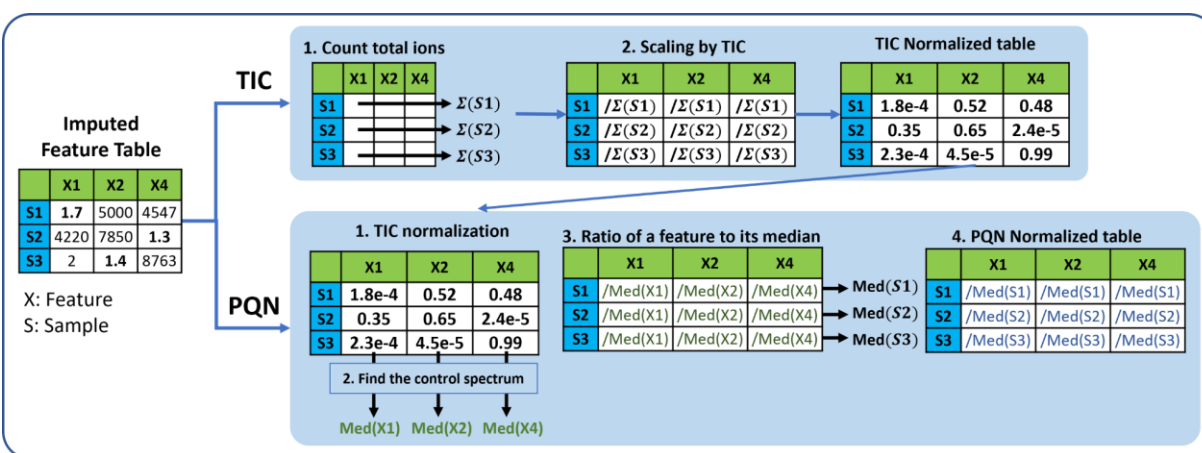
1070 1071 **Box 4 - Normalization**

Normalization of metabolomics data can rely on either chemical or mathematical strategies. The chemical method, using internal standards and quality controls, is popular in targeted analysis as it effectively balances metabolite concentrations across sample sets and batches. However, for non-targeted metabolomics, mathematical approaches are more popular^{117,120}. There are several mathematical normalization methods, each with its strengths and limitations. The selection of a normalization method depends on the specific conditions and requirements of your dataset:

1. **Unit Normalization**¹²¹ and **TIC Normalization:**
Simple and computationally efficient methods useful for large datasets. They equalize

the total sum of signal intensities across each sample. They assume that the abundance of most features does not change significantly across different samples or experimental conditions and their effectiveness decreases with large global changes in metabolite levels (e.g., due to differences in metabolite level such as healthy versus diseased, sample preparation, or instrument sensitivity). TIC normalization might over-correct disease samples with lower intensity reducing the differences between healthy and conditions.

2. **PQN**¹¹⁹: Recommended when significant size effects are present or when internal normalization disrupts relative peak information¹¹⁷. Among several LC/MS-based normalization methods, including Contrast Normalization, Cubic Splines, Cyclic Loess, PQN has been identified as the best performer in reducing sample-to-sample variations¹²⁰.
3. **Common Components and Specific Weights Analysis**¹²² (**CCSWA**): A viable alternative when QC and sample data differ.



```
norm_TIC <- normalization(imp, #performing normalization on transformed imputed data
                           method = "sum")$newXtrain
```

```
head(norm_TIC,n=3)
dim(norm_TIC)
print(paste('No.of NA values in Normalized data:',sum(is.na(norm_TIC)== T)))
```

```
norm_pqn <- normalization(imp,
                          method = "pqn")$newXtrain
```

```
head(norm_pqn,n=3)
dim(norm_pqn)
print(paste('No.of NA values in Normalized data:',sum(is.na(norm_pqn)== T)))
```

Figure 8: Normalization Techniques: Graphical representation of Total-ion-current (TIC) and Probabilistic Quotient normalization (PQN) methods, accompanied by corresponding R code snippets.

1072 3.2.5. Scaling

1073 Scaling methods in metabolomics aim to adjust the range of peak abundances between
 1074 features¹¹⁷. This is done by normalizing the intensities of each feature by a scaling factor,
 1075 effectively adjusting for fold differences between features¹²³. Additional information on scaling

1076 factors can be found in **Box 5** along with the graphical representation of scaling shown in **Figure**
1077 **9**.

1078

1079 **Step 23: Center-Scaling**

1080 We apply center-scaling to the imputed data. This allows for a consistent spread of the data,
1081 accounting for differences in offset between high and low-abundant features.

1082

1083 In R, the scale function offers different options for centering and scaling data:

1084 ● When center = TRUE, centering is achieved by subtracting the column means (excluding
1085 NAs) of the data from their respective columns (each column referring to a feature).
1086 Centering ensures that the fluctuations in the data are centered around zero instead of
1087 the mean of the metabolite concentrations¹²³.

1088 ● If center = TRUE and scale = TRUE: then scaling is performed by dividing the centered
1089 columns by their standard deviations.

1090 ● If center = FALSE and scale = TRUE: scaling is done by dividing each column by its root
1091 mean square.

1092 ● If scale = FALSE, no scaling is performed.

1093

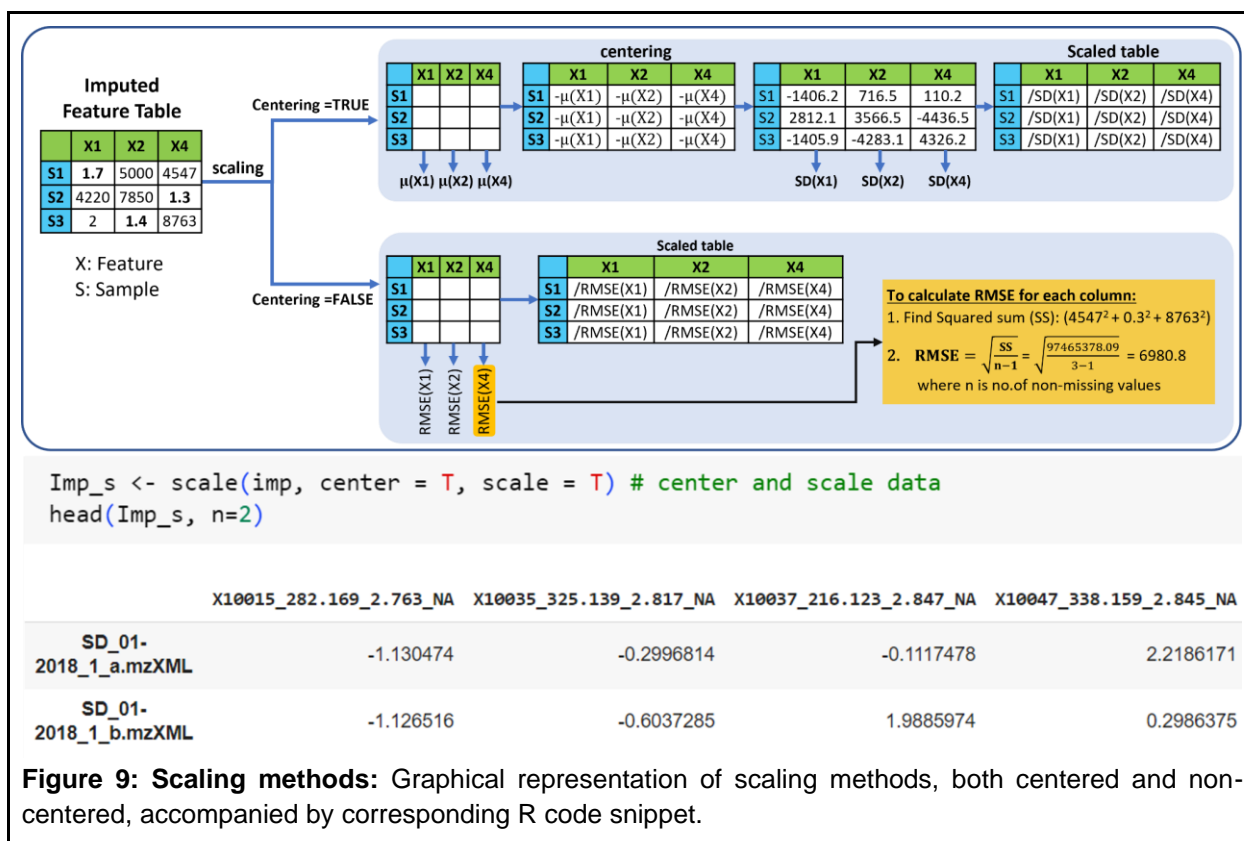
1094 ▲ **CRITICAL:** Since scaling introduces negative values, trying a PCoA with the Bray-Curtis
1095 difference on scaled data will trigger an error.

1096

1097 **Box 5 - Scaling**

Scaling methods can be categorized into two subclasses based on the scaling factor used¹²³.

1. **Using data dispersion methods**, such as standard deviation (SD), for scaling:
Examples: Autoscaling¹²⁴ and Pareto scaling¹²⁵. Autoscaling ensures equal variance (such as SD=1) for each variable, while Pareto scaling uses the square root of SD as the scaling factor.
2. **Using size measures**, such as the mean, for scaling:
Examples: Level scaling and Poisson scaling. Level scaling converts metabolite concentration changes relative to the mean concentration, while Poisson scaling scales each feature by the square root of the mean^{123,126}.



1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119

Step 24: Choosing data for further analysis

(User Input Required)

Upon executing this step, an overview table is generated, offering a list of the dataframes produced during each phase of data processing along with its respective metadata tables. This includes stages like the initial raw data (Raw Data), post-blank removal data (Blank Removed Data), post-imputation data (Imputed Data), and various normalization stages (TIC Normalized, PQN Normalized, Scaled Data).

To proceed, users must select their dataset of interest by entering the corresponding index number. The chosen dataset will be stored under the `cleaned_data` variable and the corresponding metadata will be taken under the `metadata` variable. These dataframes will be used in subsequent univariate and multivariate analytical steps. This allows the user to:

- **Explore Multiple Datasets:** Easily switch between datasets to examine the effects of different processing steps.
- **Tailor Analyses to Dataset Characteristics:**
 - TIC normalized data is apt for some univariate statistical tests, especially when analyzing the relative abundance of specific features or metabolites across samples without the comparison being skewed by samples that just have overall higher or lower intensities. Also, when using normalized data for multivariate techniques like PCA, it is important to ensure that a few dominant features do not skew the overall results.

1120 ○ Using scaled data in multivariate techniques like PCA prevents high variance
1121 features from dominating. Additionally, machine learning techniques relying on
1122 distance measures, like k-means or k-nearest neighbors, benefit from scaled data
1123 to ensure uniform feature influence.

1124 However, it is important to note:

- 1125 ● Imputation is not advised if one plans to execute a PCoA using the Jaccard distance since
1126 Jaccard transforms data into binary (0 and 1). Without zeros, it results in a table full of
1127 ones.
- 1128 ● Since scaling introduces negative values, trying a PCoA with the Bray-Curtis difference
1129 on scaled data will trigger an error.

1130
1131 For the purpose of this tutorial, we will use the `scaled_data` as our `cleaned_data` and the
1132 respective `metadata` variable is `md_Samples`. However, users are encouraged to experiment
1133 with different datasets.

1134 **3.3. Multivariate Statistics: ● Timing 50-60 mins**

1135 After data cleanup, we will use multivariate statistical analyses to allow for a deeper exploration
1136 of samples. The techniques showcased in our workflow are:

- 1137 ● **PCA and PCoA:** Principal Component Analysis (PCA) and Principal Coordinate Analysis
1138 (PCoA) are fundamental methods for discerning trends in your data. Coupled with
1139 Permutational Multivariate ANOVA (PERMANOVA), these techniques enable a
1140 comprehensive exploration of sample similarity by calculating correlations or distance
1141 matrices.
- 1142 ● **Hierarchical Clustering Analysis (HCA) and Heatmap:** This combination is ideal for
1143 hypothesis generation by providing an initial data overview. HCA builds a dendrogram
1144 representing the dataset, where individual samples are clustered based on similarity. A
1145 heatmap arranged according to the sample or feature similarities defined in the
1146 dendrogram creates a clear visual depiction of sample clusters.
- 1147 ● **Supervised Classification Techniques:** We use RF as a key supervised classification
1148 technique in this protocol. For advanced users interested in further exploration, additional
1149 instructions on XGBoost and hyperparameter tuning are provided in a separate Jupyter
1150 Notebook. The link to this additional notebook can be found in the main notebook and the
1151 file is available in our GitHub Repository. Additionally, we would like to point to Partial
1152 Least Squares - Discriminant Analysis (PLS-DA), another supervised multivariate
1153 technique that is frequently used in metabolomics studies simply due to the availability of
1154 the model in several software packages and ease of use with default settings. It handles
1155 collinear and noisy data well and offers comprehensive results such as classification
1156 prediction accuracy, scores and loadings plots. Yet, its prediction accuracy may lag behind
1157 methods like RF, especially with datasets handling fewer features. Therefore, PLS-DA
1158 might not be suitable for those who want to significantly reduce the feature numbers and
1159 then use the model on them¹²⁷. While we do not dismiss the utility of PLS-DA, we suggest

1160 considering alternative models. For a comprehensive comparison of different machine
1161 learning-based classification tools, we recommend the study of Mendez *et al.* in which
1162 they evaluate eight machine learning algorithms across ten clinical metabolomics datasets
1163 for binary classification¹²⁸.

1164 **Step 25: Installing Packages for Multivariate Analyses** ● **Timing 5-10 mins**

1165 To start our multivariate analysis, we first install and load the necessary R packages for this
1166 section: “BiocManager”, “ComplexHeatmap”, ggsci,”dendextend”, “NbClust” and “cowplot”.

1167 **3.3.1. PCoA with PERMANOVA**

1168

1169 **PCoA: Principal coordinates analysis**

1170 PCoA is a popular ordination technique used alongside PCA to visualize sample similarities by
1171 calculating distance matrices between samples. PCoA groups samples based on their
1172 dissimilarity or distances whereas PCA focusses on their correlation or covariance¹²⁹. The
1173 process begins by computing a dissimilarity matrix to capture the sample differences. This matrix
1174 is then transformed using multidimensional scaling (MDS) to produce a new set of points called
1175 Principal Coordinates (PCos) in a lower-dimensional space. The distance between samples in
1176 these coordinates reflects the original sample differences¹³⁰. It is important to mention that MDS
1177 can be categorized into metric MDS (as in PCoA) and non-metric MDS¹²⁹. In this protocol, we
1178 focus solely on metric MDS and more information can be found in **Box 6**. For graphical illustration
1179 of PCoA, see **Figure 10**.

1180 **Step 26: Prepare Data**

1181 Make sure that the metadata (`metadata``) and the feature quantification table
1182 (`cleaned_data``) are in the same order. Also, verify that the sample names (row names) in
1183 both data tables are identical and in the same order using `identical()` function. It should
1184 return TRUE.

1185

1186 **Step 27: Calculate Pairwise Distances and Perform PCoA**

- 1187 ● Calculate pairwise Euclidean distances across all samples in the feature quantification
1188 table using the `vegdist()` function from the ‘vegan’ package⁸⁴. Store the resulting
1189 distance or dissimilarity matrix as ‘dism’.
- 1190 ● Apply the `cmdscale()` function from the base R ‘stats’ package to perform MDS on the
1191 distance matrix ‘dism’, considering 10 PCos (k=10).

1192

1193 ▲ **CRITICAL:** The `vegdist()` function offers various methods such as “manhattan”, “euclidean”,
1194 “canberra”, “bray”, “jaccard”, “gower”, “binomial”, “chisq” for distance calculation. Using euclidean
1195 distance for PCoA is equivalent to performing PCA. However, using `vegdist(“euclidean”)`
1196 and `cmdscale()` cannot provide loadings information. For a comprehensive PCA with both
1197 loadings and scores, use the `prcomp()` function such as `pca_result <-`

1198 `prcomp(cleaned_data, center = FALSE, scale. = FALSE)`. Since our ``cleaned_data``
1199 we use is already centered and scaled, we can set these parameters to FALSE. For loadings and
1200 PC scores, you can access ``pca_result$rotation`` and ``pca_result$x`` respectively.

1201

1202 **Step 28: Analyze PCoA Results**

1203 Examine the list generated by the `cmdscale()` function, which includes the following elements:

- 1204 • `'points'` (`Pcoa$points`) represents the data matrix with the given PCos
- 1205 • `'eig'` (`Pcoa$eig`) indicates the eigenvalues computed for the PCos, which describe the
1206 variance explained by each PCo.

1207

1208 **Step 29: Plot PCoA Scores**

1209 *(User Input Required)*

1210

1211 Using the ``ggplot2`` library, create a PCoA Scores Plot. Here, the samples are color-coded
1212 based on the `'ATTRIBUTE_Month'` attribute. To view the sample distribution of different
1213 attributes, simply adjust the line: `interested_attribute_pcoa = 'ATTRIBUTE_Month'`.
1214 Importantly, the aspect ratio of the plot's axes is maintained to ensure accurate representation,
1215 in line with recommendations by Nguyen and Holmes¹³¹.

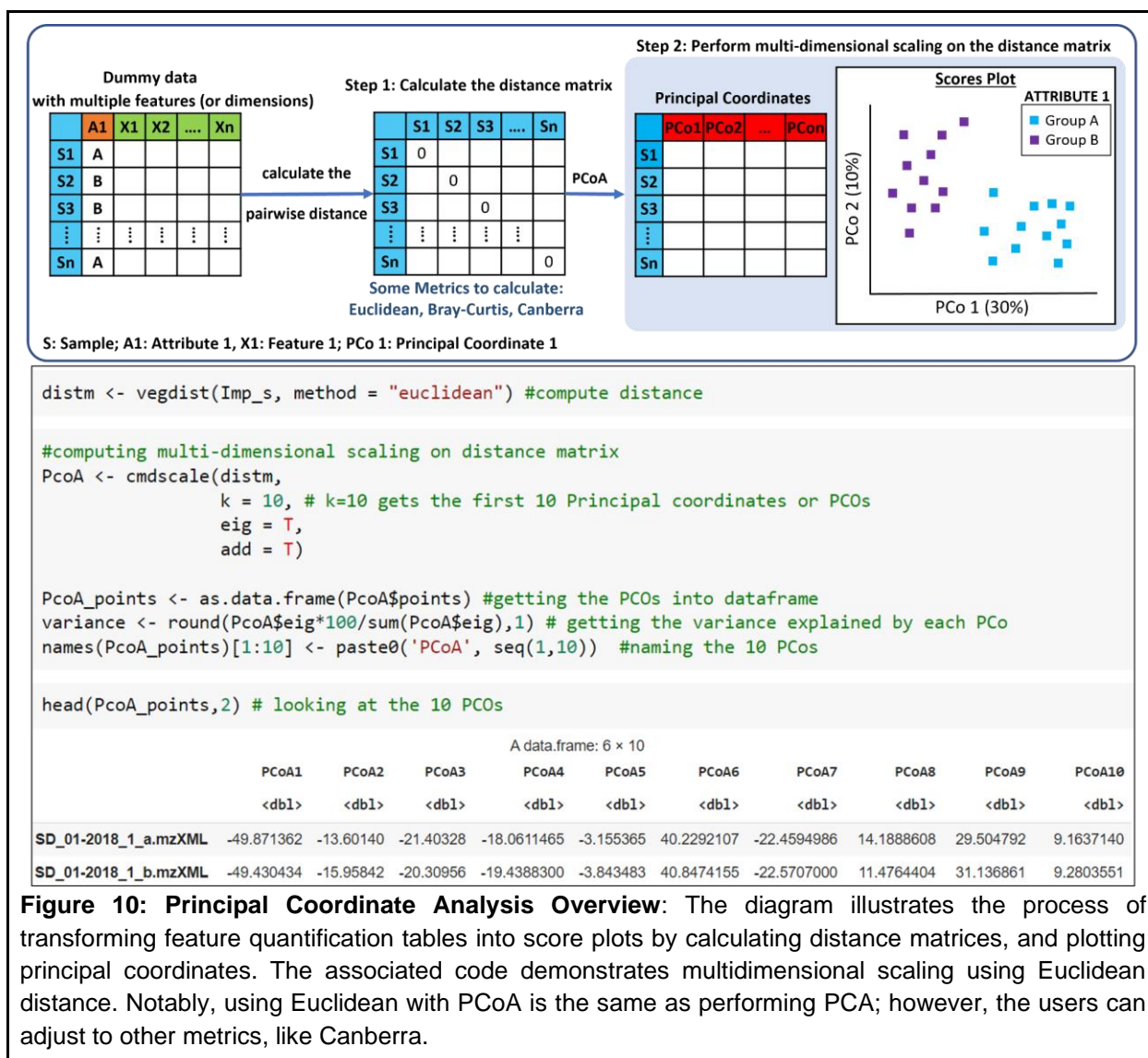
1216

1217 **Box 6 - Principal Coordinate Analysis (PCoA)**

PCoA offers an advantage over PCA by allowing various distance metrics beyond the Euclidean distance. This flexibility provides different insights into the data pattern based on the chosen dissimilarity measure. For example, when working with categorical data and sparse matrices containing numerous zeros, distance metrics such as Hamming distance and Jaccard distance outperform the Euclidean distance^{130,132,133}. Akin to phylogenetic distance measures such as UniFrac distance¹³⁴ used in the microbial ecology field, chemical distance matrices are emerging that make use of cosine MS/MS similarity between features¹³⁵ or chemical similarity derived from CSI:FingerID¹³⁶.

While PCoA effectively reveals chemical trends among samples by working with different distance matrices, it cannot provide direct information about the relationship between features and principal coordinates, unlike PCA which offers 'loadings' information¹³⁷. To discern associated features in such contexts, it is recommended to complement PCoA with other methods like a Heatmap overview, Random Forest analysis or any of the univariate techniques discussed in this protocol.

In addition, to assess the impact of a specific feature on the dispersion of samples along a particular PCoA axis, an indirect analysis can be performed. This involves correlating or regressing the PCoA values of the samples with the corresponding sample scores of the variable of interest¹³⁸. For instance, in our case, to evaluate the influence of Feature 1 on PCo1, we can create a scatter plot by plotting the original values of Feature 1 (sample scores) for all samples against the PCo1 values for all samples. The points on the plot can be colored based on the sampling period. By examining any trends or correlations in the plot, we can observe how the diversity of samples changed during the sampling period.



1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229

PERMANOVA: Permutational multivariate ANOVA

In multivariate analysis like PCA, it is crucial to measure confidence in observed relationships or separation between objects. This is often achieved via statistical significance tests, which provide a p-value as a measure of the confidence level. For ordination techniques that do not assume a specific data distribution, parametric statistical testing is not applicable¹³⁹. In such cases, resampling methods such as bootstrap, jackknife¹⁴⁰, and permutation tests¹⁴¹ are used to assess the statistical confidence of the results. These methods generate multiple samples or permutations from original data to estimate variability and assess the significance of observed relationships¹³⁹.

1230 Alternatively, non-parametric methods such as PERMANOVA (Permutational Multivariate
 1231 Analysis of Variance) can be used¹⁴². PERMANOVA allows for multivariate ANOVA and tests for

1232 differences between object classes. It enables any dissimilarity metric and calculates a test
1233 statistic by comparing the dissimilarities between objects within and between classes. Here, the
1234 p-values are determined through permutation¹³⁹.

1235 **Step 30: Testing for Homoscedasticity**

1236 *(User Input Required)*

1237 Before performing PERMANOVA, it is important to validate the homogeneity of group
1238 dispersions, often termed as 'Homoscedasticity'. This test ensures that each group exhibits
1239 approximately equal variability. Violation of this assumption might inflate the risk of Type I errors
1240 (false positives).

1241 If the group dispersions are homogenous, you can proceed with PERMANOVA with greater
1242 confidence. However, disparate dispersions require a more cautious interpretation of
1243 PERMANOVA results, given their higher susceptibility to Type I errors. In such cases, exploring
1244 alternative distance measures, data transformations, or delving into potential biological reasons
1245 for the dispersion differences might offer a more comprehensive analysis. To know more about
1246 multivariate dispersions, see **Box 7**. For a visual representation of assessing multivariate
1247 dispersion and conducting the PERMANOVA analysis in R, refer to **Figure 11**.

1248
1249 Procedure to Evaluate Homoscedasticity:

- 1250 ● As a first step, the user needs to specify the attribute group for assessing group
1251 dispersions. Since we are looking for group dispersions, it is important to select a
1252 categorical metadata column (for example, 'ATTRIBUTE_Month') and avoid choosing
1253 continuous attributes, such as 'ATTRIBUTE_Injection_order'.
- 1254 ● Similar to Step 27, we compute a distance matrix ('dism') using the feature quantification
1255 table and the selected attribute. For simplicity, we use the Euclidean distance in this
1256 instance.
- 1257 ● Using the `betadisper()` function from the `vegan` package, we evaluate group
1258 dispersion against the chosen attribute group.
- 1259 ● The dispersion model is then visualized to offer a clearer perspective.
- 1260 ● Lastly, an ANOVA is executed on the dispersion model. A significant p-value ($P < 0.05$)
1261 indicates a violation of the PERMANOVA's foundational assumptions. Conversely, a
1262 non-significant result suggests that PERMANOVA is a suitable choice for the given
1263 attribute.
- 1264 ● The resulting p-value for 'ATTRIBUTE_Month' is significant, indicating the presence of
1265 group dispersions among different months. This violates the PERMANOVA assumption.
1266 When PERMANOVA is performed for this attribute, the PERMANOVA results require a
1267 more cautious interpretation.

1268 1269 **Step 31: Conduct PERMANOVA Test**

- 1270 ● Use the `adonis2()` function from the `vegan` package⁸⁴ to conduct a PERMANOVA test.
1271 The `adonis2` function allows for the analysis and partitioning of sums of squares using
1272 dissimilarity measures.

1273
1274
1275
1276
1277
1278
1279

- Apply the `adonis2` function on the dissimilarity matrix ('dism') and the previously chosen metadata column 'ATTRIBUTE_Month'. This helps in investigating if there are significant differences among the samples collected during three different months.
- Interpret the resulting p-value. In our case, we obtained a p-value of 0.001, indicating a significant difference between the samples.

Box 7 - Dispersion Analysis

In the case of balanced sample sizes across groups, PERMANOVA identifies differences in group centroids, thus reflecting shifts in the multivariate distribution of sample units within the chosen resemblance space. Hence, the type of dissimilarity measure you choose is crucial. For example, unlike Euclidean distance, measures like Jaccard or Bray-Curtis highlight the similarity in species composition and do not focus on the central tendency such as the mean-variance relationship. On the other hand, PERMDISP is specifically tailored to detect variations in multivariate dispersions. Therefore, when analyzing your data, use PERMANOVA to understand group centroid shifts and PERMDISP to evaluate dispersion differences¹⁴³.

a)

```
dispersion_model <- betadisper(dism, group)
disp <- anova(dispersion_model)
disp["significant"] <- ifelse(disp$`Pr(>F)`<0.05, "Significant", "Non-significant")
disp
```

A anova: 2 × 6

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	significant
	<int>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>
Groups	2	37529.38	18764.6917	31.56342	1.890734e-12	Significant
Residuals	177	105227.82	594.5075	NA	NA	NA

b)

```
adonres <- adonis2(dism ~ group)
adonres
```

A anova.cca: 3 × 5

	Df	SumOfSqs	R2	F	Pr(>F)
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
group	2	385128.9	0.236643	27.43527	0.001
Residual	177	1242339.1	0.763357	NA	NA
Total	179	1627468.0	1.000000	NA	NA

Figure 11: Multivariate Dispersion and PERMANOVA Analysis: a) Code snippet for testing multivariate dispersions within the 'group', specifically referencing the 'ATTRIBUTE_Month' column (Dec, Jan, Oct) from the metadata. b) Code snippet for executing PERMANOVA to analyze variations between the aforementioned groups.

1280 **Step 32: Define a Function for Streamlined Analysis**

1281 To facilitate quicker analysis and avoid rewriting from step 27 to step 31 for testing different
1282 parameters, we defined a function, `plotPCoA()`. This function performs a principal coordinates
1283 analysis (PCoA) using a chosen distance metric, calculates a PERMANOVA, and plots the results
1284 in a 2-D graph. Additionally, it assesses group dispersion prior to the PERMANOVA calculation
1285 and displays the significance result in the resulting plot as well.

1286
1287 The function has the following parameters:

- 1288 - `ft` refers to the desired feature quantification table.
- 1289 - `md` refers to the respective metadata.
- 1290 - `distmetric` is the distance metric of choice.
- 1291 - `category_permanova` is the desired metadata group for PERMANOVA calculation.
- 1292 - `pcoa_category_type` indicates whether the group type is categorical or continuous.
- 1293 - `category_pcoa_colors` specifies the metadata attribute for coloring the samples.
- 1294 - `cols` are the desired colors for the groups.
- 1295 - `title` is the title of the plot.

1296
1297 Additionally, we have created another simple custom function `save_as_svg()`, to store plots in
1298 SVG format utilizing the `svglite` function. This custom function can be used as
1299 `save_as_svg(filename, desired_plot, plot_width, plot_height,`
1300 `plot_background)`. Throughout the notebook, you will observe this function being employed
1301 post each plot creation to save the visualizations.

1302

1303 **Step 33: Applying `plotPCoA()` function on different dataframes**

1304 *(User Input Required)*

1305 In this step, the user can specify the variables as mentioned in the previous step. Here is an
1306 example of how to use the `plotPCoA()` function:

```
plotPCoA(  
  ft = cleaned_data,  
  md = metadata,  
  distmetric = "euclidean",  
  category_permanova = "ATTRIBUTE_Month",  
  pcoa_category_type = 'categorical',  
  category_pcoa_colors = "ATTRIBUTE_Month",  
  cols = c('orange', 'darkgreen', 'red', 'blue', 'black'),  
  title = 'Principal coordinates plot')
```

1307 **Step 34: Get PCoA plots after each data cleanup step**

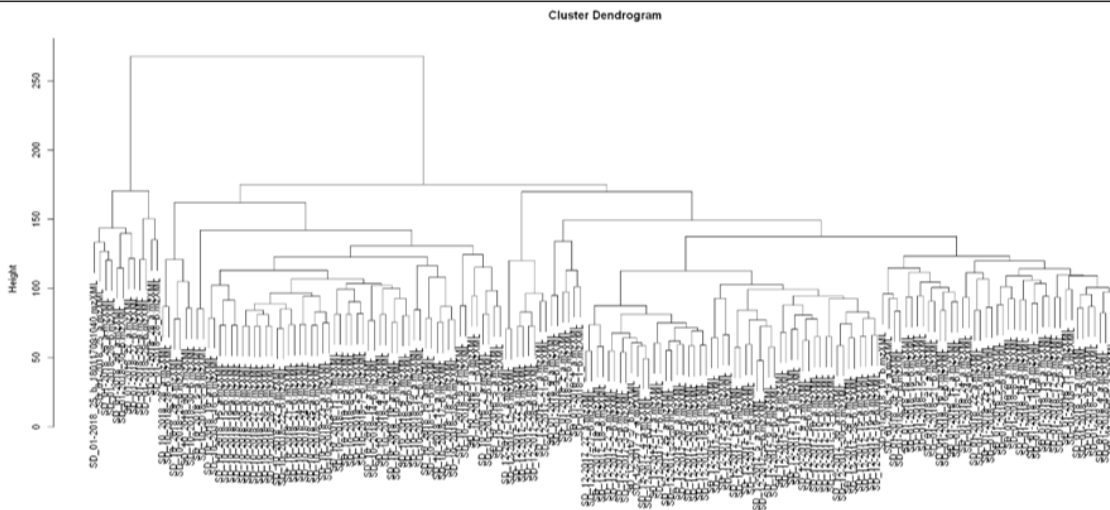
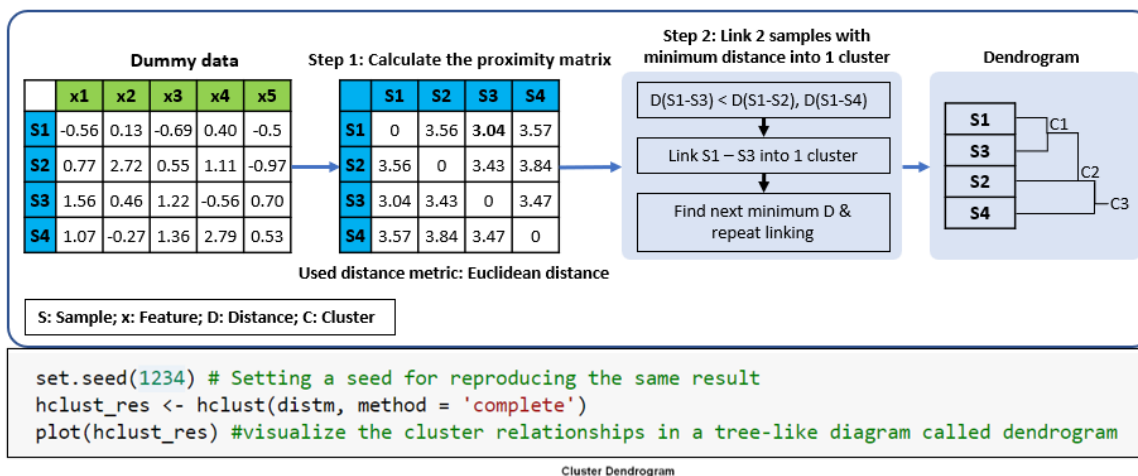
1308 *(User Input Required)*

1309
1310 In this step, the user can specify parameters such as the distance metric, attribute for
1311 PERMANOVA calculation, attribute to color the PCoA scores, the category of the chosen attribute,

1312 similar to the previous plotPCoA step. These inputs will be taken to produce an overview of PCoA
1313 plots for all steps of data cleanup.

1314 3.3.2. Hierarchical Cluster Analysis

1315 Clustering, such as Hierarchical Cluster Analysis (HCA), is an unsupervised classification method
1316 commonly used in metabolomics to determine the similarity between samples based on their
1317 chromatograms or other characteristics. Unlike PCA, which focuses on capturing the maximum
1318 variance between samples, clustering aims to group samples with “similar” profiles. The results
1319 are often visualized as dendrograms¹³⁰ as shown in **Figure 12**.



1320 **Figure 12: Dendrogram Generation and Analysis:** The figure illustrates a dendrogram, as a result of
1321 applying HCA to a feature quantification table (e.g., ‘cleaned_data’). From this data, a proximity matrix (or
1322 the distance matrix) is calculated (see steps 27 and 30), which subsequently guides the dendrogram
1323 creation. Accompanying the illustration is the related code for the cluster generation and dendrogram
1324 visualization. The distance matrix ‘distm’ is calculated via Euclidean distance in step 30, though alternative
1325 metrics can be chosen by the user. The resultant dendrogram is displayed, initially partitioning samples into
1326 two primary clusters: a smaller cluster from a subset of samples (corresponding to samples from January
1327 in our example data) and a larger subsequent cluster. Distinct sub-clusters within these main clusters are
1328 also discernible.
1329

1330 **Step 35: Setting the Plot Size**

1331 First, we need to define the size of the output plot, as dendrograms are typically larger in size.
1332 Adjust the plot size accordingly to ensure a clear and comprehensible visualization.

1333 **Step 36: Executing HCA**

1334 Next, we use the `hclust()` function from the 'stats' package to perform HCA. The function is
1335 applied to the distance matrix 'dism', calculated based on the feature quantification table
1336 ('cleaned_data') using a specified distance metric (e.g., Euclidean, Canberra). The 'method'
1337 argument in `hclust()` denotes the linkage method used for measuring the distance between
1338 clusters (e.g., complete, single, average). We use the default 'complete' method, which calculates
1339 the maximum distance between clusters before combining them. Once HCA is completed, a
1340 dendrogram is generated. This dendrogram shows split or merge distances as 'height' along the
1341 y-axis, providing a visual representation of the cluster formation.

1342

1343

1344 **Step 37: Cutting the Dendrogram**

1345 *(User Input - Optional)*

1346 Similar to k-means clustering, which seeks to establish k clusters with minimum within-cluster
1347 variation, we can cut the dendrogram into a specified number of clusters using the `cutree()`
1348 function. However, we need to initialize the clustering with random k clusters. For our sample
1349 dataset, we define `k=4` with the `cutree()` function, to create four clusters. The user can change
1350 the number of clusters.

1351 **Step 38: Coloring the Dendrogram**

1352 Finally, we can extract the cluster allocation information and color the dendrogram according to
1353 the clusters. For our data, the dendrogram suggests two main splits, resulting in four distinct
1354 clusters.

1355 **Step 39: Determining the Optimal Number of Clusters**

1356 Here, we use heuristic methods similar to those applied in k-means clustering to determine the
1357 optimal number of clusters. For this purpose, we use the Elbow approach and average silhouette
1358 method using the `fviz_nbclust()` function from the 'factoextra' package.

1359

- 1360 ● The Elbow method calculates the total within-cluster sum of squares (WSS) for an
1361 increasing number of clusters. WSS signifies the sum of distances between data points
1362 and their corresponding centroids within each cluster. Lower WSS values indicate within-
1363 cluster variation.
- 1364 ● The resulting Elbow plot presents the WSS on the y-axis and the number of clusters on
1365 the x-axis. Lower WSS values suggest minimum within-cluster variation and better
1366 clustering. However, the 'elbow' point is considered as an indicator of the optimal number

1367 of clusters, as further cluster additions do not significantly improve the clustering or
1368 decrease the WSS. For our example data, this method suggests 3 or 4 clusters. However,
1369 defining the 'elbow' can be subjective.

- 1370 • An alternate approach is the average silhouette method, which assesses clustering quality
1371 by determining how well each data point fits within its assigned cluster. In our case, this
1372 method proposes two primary clusters.

1373
1374 Both the Elbow and Silhouette methods provide global insights without learning from the data,
1375 given their unsupervised nature. But, there are more sophisticated techniques like the gap-
1376 statistic which refines the heuristic concepts behind the Elbow and Silhouette techniques and
1377 uses a statistical procedure to estimate the optimal cluster count¹⁴⁴. However, all these methods
1378 serve as guidelines rather than definitive answers. In practice, users might choose cluster
1379 numbers based on context, for example, in our case with seven sample areas, opting for seven
1380 clusters can be insightful. Later, one can check whether these clusters correspond to known
1381 sample groups.

1382 3.3.3. Heatmaps

1383 Heatmaps are generally used to visualize complex data or discern patterns across a high-
1384 dimensional dataset. They are commonly used in bioinformatics¹⁴⁵, particularly in gene expression
1385 analysis and visualizing genomic datasets, owing to their ability to effectively represent thousands
1386 of data points¹⁴⁶. This makes them equally suitable for mass spectrometry-based metabolomic
1387 experiments. Heatmaps are efficient in pattern recognition due to their color-coded matrix
1388 elements and adjacent dendrograms, which indicate functional relationships between variables
1389 and samples¹⁴⁷. For more information on heatmap, see **Box 8**. To see the resulting heatmap
1390 generated by the R code in the Notebook, refer to **Figure 13**. In this section, we will show how to
1391 incorporate hierarchical clustering into our heatmap.

1392

1393 **Step 40: Preparing Metadata for Heatmap**

1394 *(User Input Required)*

1395 To start with, determine which metadata columns or attributes will be used to decorate the
1396 heatmap. In our case, we specified the following attributes: 'ATTRIBUTE_Year',
1397 'ATTRIBUTE_Month', and 'ATTRIBUTE_Sample_Area'. The user can select any number of
1398 attribute columns from their metadata as they see fit for the heatmap. A new dataframe is created
1399 comprising the chosen metadata.

1400 **Step 41: Generating annotations for Heatmap**

1401 *(User Input - Optional)*

1402 For distinct visualization, this step assigns unique colors to each category within chosen attributes
1403 from the previous step. We have created a function `generate_colors()`, which utilizes a
1404 predefined color-blind-friendly palette of 10 colors to assign colors to these unique groups. Users

1405 can modify these colors if desired. After assigning colors to the subset dataframe, we use this
1406 information to decorate the heatmap with annotations from the `HeatmapAnnotation()` function
1407 in the 'ComplexHeatmap' package.

1408 **Step 42: Creating the Heatmap**

1409 *(User Input - Optional)*

1410 To create the heatmap, apply the Heatmap function from the ComplexHeatmap package on the
1411 transposed ``cleaned_data`` (as previously chosen in step 24). This arranges the features in rows
1412 and samples in columns.

- 1413 • For the heatmap, the color intensity represents the feature intensities, with the intensity
1414 scale ranging from 0 (blue) to 1 (dark red), and 0.5 represented as white. This color coding
1415 allows for a visual comparison of feature intensity variations across samples.
- 1416 • The clustering on the y-axis is based on Euclidean distance
1417 (`clustering_distance_rows = "euclidean"`, `clustering_distance_columns =`
1418 `"euclidean"`). However, other distance measures such as Manhattan, Minkowski,
1419 Canberra, or even Jaccard for binary data, can be chosen based on specific needs.
- 1420 • The 'complete' linkage method is used for clustering (`clustering_method_rows =`
1421 `"complete"`, `clustering_method_columns = "complete"`).

1422

1423 **Step 43: Refining Data Clustering with k-means**

1424 Further refine data clustering by incorporating the built-in k-means function within the heatmap as
1425 parameters for row and column clustering (`row_km = 5`, `column_km = 4`). To ensure robustness,
1426 perform multiple repeats (`row_km_repeats = 100`, `column_km_repeats = 100`).

1427 **Step 44: Extracting Features from Each Cluster**

1428 With a higher number of features, it is difficult to interpret the clustering or labeling of features on
1429 the heatmap. To address this, extract the features from each cluster into a separate dataframe.
1430 This dataframe containing combined feature names (``XFeatureID_m/z_RT_GNPS_annotations``)
1431 and their respective cluster assignments can be saved as a CSV file for further interpretation. For
1432 example, one could merge these cluster assignments with the feature quantification table for
1433 import into Cytoscape along with the FBMN and use these cluster assignments for coloring slices
1434 in node pie charts.

1435 **Box 8 - Heatmap**

Although widely used, traditional cluster heatmaps also have limitations. Their data representation in two-dimensional format can be restrictive when processing complex multidimensional data. Furthermore, their static nature does not allow for data to be sorted along different axes, filtered, or focused on specific elements, making the representation of a vast

number of elements quite challenging. Regardless of these limitations, heatmaps are preferred in biological and biomedical data representation because their visual format simplifies data interpretation and comparison. To overcome these limitations, more advanced versions such as XCMS interactive heatmaps are available that offer a more versatile and dynamic data visualization experience¹⁴⁷.

```
# set the parameters for the type of clustering to perform. You can play with different options
set.seed(1234)
hmap <- Heatmap(
  t(cleaned_data),
  heatmap_legend_param = list(title = "Scaled/centered \n intensity"),
  col = circlize::colorRamp2(c(0, 0.5, 1),
    colors = c("blue", "white", "darkred")),
  show_row_names = FALSE, show_column_names = FALSE,
  cluster_rows = TRUE, cluster_columns = TRUE,
  show_column_dend = TRUE, show_row_dend = TRUE,
  row_dend_reorder = TRUE, column_dend_reorder = TRUE,
  clustering_distance_rows = "euclidean", # you can change the distance here
  clustering_distance_columns = "euclidean", # you can change the distance here
  clustering_method_rows = "complete",
  clustering_method_columns = "complete",
  width = unit(100, "mm"),
  top_annotation = colAnn)

ComplexHeatmap::draw(hmap, heatmap_legend_side="right", annotation_legend_side="right")
```

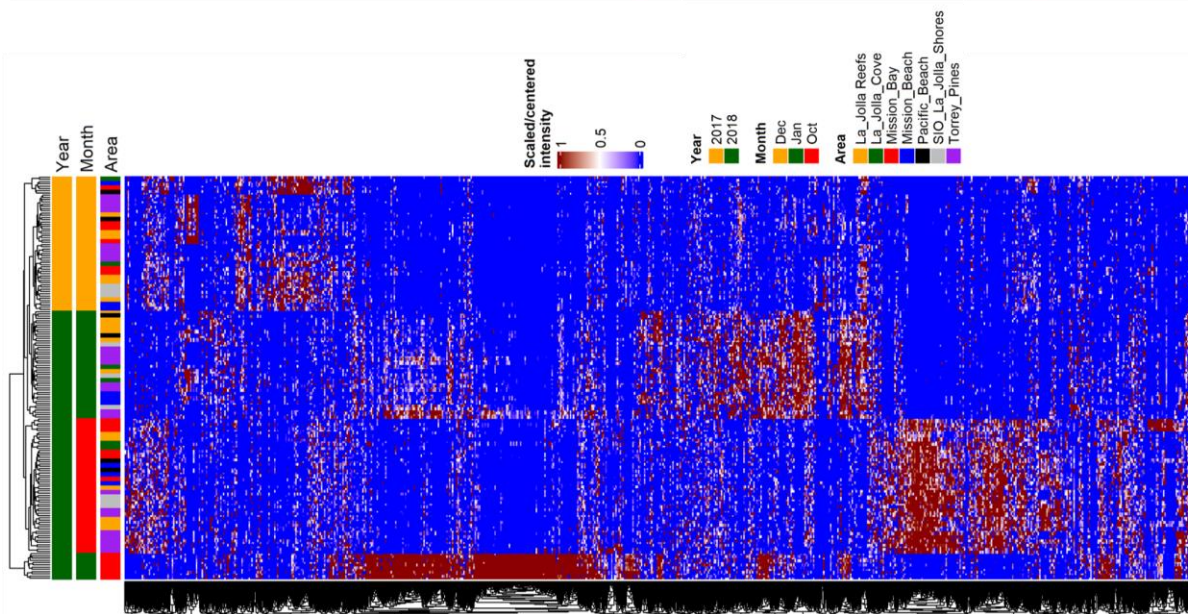


Figure 13: Heatmap Visualization and Construction: This figure presents both the R code snippet used for heatmap creation and the resultant heatmap itself. To facilitate a comprehensive view, the heatmap is oriented horizontally. The feature quantification table used here is the scaled table and feature intensities are color-coded, ranging from blue (0) to red (1). Annotations at the heatmap's top delineate clustering based on variables like year, month, and sample area.

1437 3.3.4. Supervised Classification: Random Forest

1438
1439 Unsupervised analysis allows for the discovery of groups or trends in the data without prior
1440 assumptions about any predetermined labels or categories, whereas supervised analysis involves
1441 the use of labeled data to guide the analysis toward specific objectives such as biomarker
1442 discovery, classification, and prediction. In supervised analysis, the algorithm is trained on labeled
1443 data to predict the response variable (or dependent variable) based on the predictor variables (or
1444 independent variables)¹⁴⁸.

1445
1446 Supervised learning is categorized into classification and regression problems based on the type
1447 of response variable: classification for categorical or discrete variables (e.g., cancer vs non-
1448 cancer samples), and regression for continuous variables. Popular supervised models in
1449 metabolomics include logistic regression, partial least square discriminant analysis (PLS-DA),
1450 support vector machines (SVM), k-nearest neighbor (KNN), and random forest (RF). Here, we
1451 focused on RF, which offers advantages such as the low risk of overfitting, ease of
1452 implementation, interpretability, and minimal hyperparameter tuning requirements¹⁴⁹. For a
1453 detailed overview of random forest, consult **Box 9**.

1454
1455 In our example provided in the notebook, we tried to classify surface seawater samples based on
1456 their different sampling sites using random forest. Here, the feature quantification table without
1457 metadata is the predictor variable, and the metadata group “Sampling Site” is the response
1458 variable. **Figure 14** provides a visualization of the Random Forest algorithm and its
1459 implementation in R.

1460
1461 **Step 45: Prepare the data for Random Forest**
1462

- First, load the ``rfPermute`` package.
- Start by merging the feature quantification table (in our example, ``Imp_s`` is chosen as
1463 the ``cleaned_data`` variable) and the corresponding metadata (``md_Samples``) into a
1464 dataframe named ``cleaned_data_with_md``. This step ensures that the samples are
1465 correctly aligned with their corresponding attributes in the metadata, which is essential for
1466 the subsequent analyses.

1467
1468
1469 **Step 46: Select the Classification Attribute for Random Forest**
1470 *(User Input Required)*
1471 Prepare the dataset used for Random Forest classification so it only contains feature intensity
1472 information and attribute of interest for classification. Here, we are classifying the samples
1473 according to different sample areas (``ATTRIBUTE_Sample_Area``). So in this step, the user is
1474 prompted to input the index number of the interested attribute to use for the classification.

1475
1476 **Step 47: Balance sample sizes**
1477 If the sample size varies among the groups, balance the size using the `balancedSampleSize()`
1478 function.

1479

1480 **Step 48: Run Random Forest**

1481 Initiate the Random Forest analysis by setting the number of trees and permutations. In our case,
1482 we used 500 trees (`ntree``) and 500 permutations (`num.rep``). Here, the primary parameters
1483 for Random Forest include the feature quantification table (without the classification data),
1484 predictor variable, balanced sample size (`sampsiz``), and tree and repetition quantities
1485 (`ntree`` and `num.rep``).

1486 With the `rfpermute`` function, there is no need for the conventional train-test split, such as the
1487 70-30 or 80-20 ratio. This user-friendly package minimizes the need for parameter tuning.
1488 Classification rates in Random Forest rely on out-of-bag (OOB) samples, which are not part of
1489 the tree-building process. This eliminates the need to split the dataset into test and train portions,
1490 maximizing the amount of information the model has to build a classifier. However, classes with
1491 unequal sample sizes, will produce a model that will tend to perform better on the larger class. To
1492 alleviate this bias, create a balanced model where the classes are represented by an equal
1493 number of samples in each tree and sampling is done without replacement using the
1494 `balancedSampsiz`` function.

1495 **▲ CRITICAL:**

- 1496 • Increasing the number of trees and permutations generally enhances the model's
1497 performance but also escalates computational costs. It is advised to start with a
1498 reasonable number of trees (e.g., 500-1000) and `num.rep`` (500-10000), then adjust
1499 based on performance.
- 1500 • When working with large data sets, R may run out of internal memory trying to perform the
1501 random forest. To work around this, adding the "`as.factor``" in the predictor variable (y),
1502 even if the class is already a factor, will alleviate the memory error.

1504 **Step 49: Evaluate model performance**

1505 After getting the RF model, we need to evaluate the model's performance using several metrics
1506 such as model accuracy, the confusion matrix, trace plot, and check for potential overfitting by
1507 comparing testing versus training accuracies.

- 1508 • The confusion matrix is the most basic summary of a Random Forest. The matrix consists
1509 of the 'original class' in rows and the 'predicted class' in columns. The diagonals represent
1510 the number of samples correctly classified in each class. The matrix also has columns that
1511 show the percent of samples that were correctly classified in a class, along with upper and
1512 lower 95% confidence intervals.
- 1513 • The trace plot shows the OOB (out-of-bag) changes as trees were added to the forest.
1514 The model should have enough trees in it so the error rate is stable. If the error rate level
1515 increases as the number of trees increases, it may be an indication of overfitting.

1517 **Step 50: Interpreting RF Results**

1518 Beyond these, the RF results can be interpreted in various ways:

- 1519 • One could plot the most impactful predictors in the model using violin plots. Here, we show
1520 the top 9 predictors
- 1521 • Compare class predictions versus the actual group in a proximity plot

- 1522
- 1523
- 1524
- 1525
- 1526
- 1527
- 1528
- 1529
- Rank features by importance using the 'Mean Decrease Accuracy' metric. This metric helps identify features whose removal significantly impacts the model's accuracy, thus marking their importance. If a feature's removal does not affect accuracy, it may be deemed less important. Features with a 'MeanDecreaseAccuracy.pval' < .05 are considered significant, implying that their absence would affect the model's performance significantly. This ranked list can also be exported as a CSV file for further analysis.

Box 9 - Random Forest

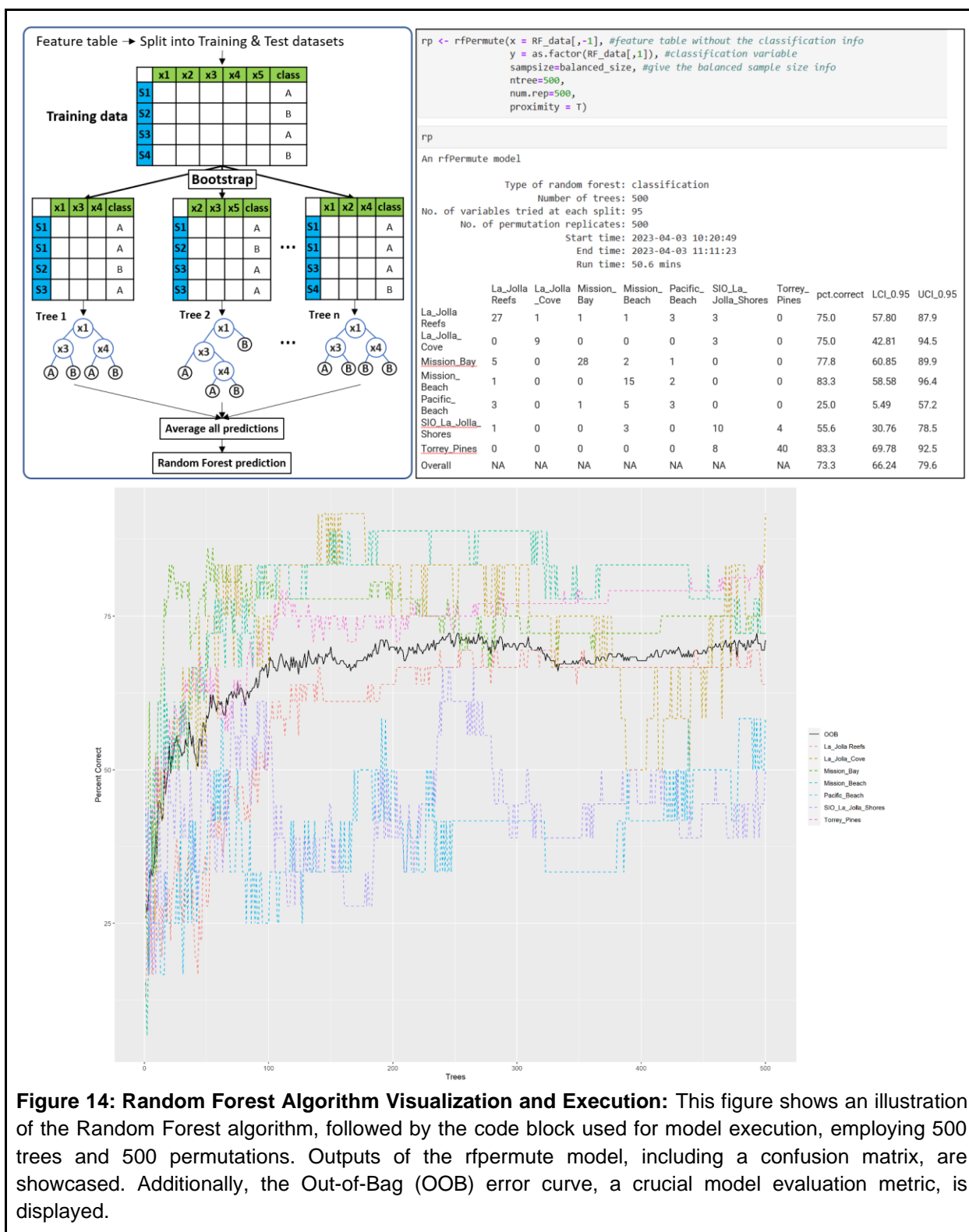
Random Forest (RF) is a powerful machine learning algorithm that operates by dividing data into fractions, building randomized tree predictors on each fraction, and aggregating these predictors together. Generally, RF uses out-of-bag (OOB) error as an estimate of the overall generalization error and obtains variable importance scores through permutation¹⁵⁰.

A unique feature of the RF algorithm is its use of OOB samples, which are the samples not used in the bootstrap sample for a particular tree. Each tree is trained on about two-thirds of the total dataset, with the remaining one-third serving as the OOB samples. The OOB error rate is a measure of prediction accuracy and helps to improve the performance of weak or unstable learners in the model¹⁵¹.

In RF, variable importance scores are obtained by permuting the values of each variable 'm' within the OOB samples and the tree is used to make predictions on these permuted OOB samples. This essentially disrupts any relationship that variable 'm' might have with the target variable. The model then compares the prediction accuracy on the variable-m-permuted OOB samples to predict accuracy on the original (untouched) OOB samples. The average of the difference in accuracy (between permuted and original OOB) across all trees in the forest gives the raw importance score for variable "m". This raw importance score is often an average value over all trees. To determine if this importance score of variable "m" is statistically significant, a z-score can be calculated by dividing the raw score by its standard error¹⁵².

In RF, there are two common metrics of variable importance used to rank features based on their predictive power: Mean Decrease Accuracy (MDA) and Mean Decrease Gini (MDG). MDA measures the decrease in model accuracy when a particular variable's values are permuted. A large decrease indicates high variable importance; MDG measures how each variable contributes to the homogeneity of the nodes and leaves in the resulting RF. A higher MDG value indicates that splitting the dataset by this variable results in purer nodes. Here, Variable Importance Projection (VIP) could be obtained by normalizing MDA, so they sum to 100, making them more interpretable on a relative scale¹⁵³.

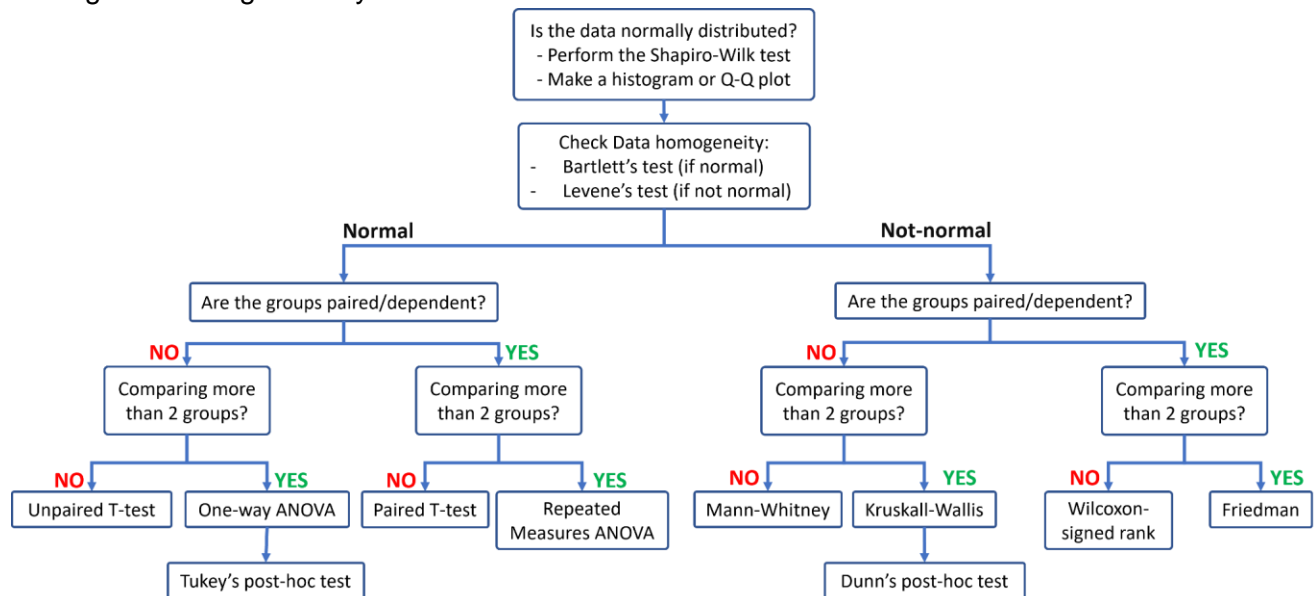
Some of the other important parameters to keep in mind to evaluate the performance of the RF model are: model accuracy, confusion matrix (a matrix showing true vs predicted class labels), trace plot, and check for overfitting by comparing testing vs training accuracy. However, supervised models may not be suitable for all data sets, especially those with few observations or unclear class distinctions. Confounding variables, related to both the predictor and response variable, can also make these models unsuitable. For instance, age and gender in a drug study can be confounding variables, leading to erroneous results if not controlled for. In such cases, using supervised models for analysis may not be appropriate.



1531 **3.4. Univariate Statistics: ● Timing 50-60 mins**

1532 While multivariate analyses offer a comprehensive overview of the data, univariate statistical
1533 analyses allow us to focus on specific attributes. Primarily, univariate analysis in metabolomics
1534 helps identify individual metabolites that significantly differ between experimental groups,
1535 potentially serving as biomarkers for certain conditions or indicators of specific biological
1536 processes. It can also reveal impacts on specific metabolic pathways if related metabolites
1537 change significantly. However, it is worth noting that univariate analysis does not account for
1538 metabolite correlations and interactions, hence, it's best used in conjunction with multivariate
1539 analysis for a holistic data interpretation.

1540 For example, our test dataset consists of numerous features collected at seven diverse sample
1541 sites. Here, univariate analyses can assess feature differences across these sites. In the case of
1542 two site comparisons, the t-test can be used to examine significant feature differences (p value <
1543 0.05). For a comparison involving more than two sample groups, we utilize ANOVA. **Figure 15**
1544 provides a flowchart that guides the selection of appropriate statistical tests based on data
1545 normality and homogeneity. In the event of significant differences, we represent these findings
1546 through a bar graph that captures the distribution of a 'significant' feature across sample
1547 conditions. Post-hoc tests are also introduced as supplementary tools to identify which groups'
1548 average values significantly differ.



1549 **Figure 15:** Flowchart detailing the selection of statistical tests for univariate analysis, based on data
1550 normality and homogeneity.
1551

1552
1553 When conducting multiple univariate tests simultaneously, as is common in metabolomics, there
1554 is an increased risk of false positives. To manage this, the False Discovery Rate (FDR) gauges
1555 the expected false positives among significant results. While the classical Bonferroni correction
1556 addresses false positives, it could increase the false negative rate. The following are some
1557 advanced methods that focus on maximizing true discoveries without escalating the false positive
1558 error rate¹⁵⁴.

- 1559 • **Benjamini-Hochberg (BH):** Commonly used in metabolomics for being less conservative
1560 than Bonferroni. It ranks p-values and adjusts them, targeting the expected false positives
1561 among all positives, rather than across all tests. It calculates FDR as Expected (False
1562 Positive/ (False Positive+True Positive)).
- 1563 • **Benjamini-Yekutieli (BY):** An iteration of BH that is suitable when tests have
1564 dependencies.
- 1565 • **Storey's q-value:** This approach estimates the proportion of true null hypotheses (i.e., no
1566 effect) among all hypotheses and then computes a q-value for each test, which is the FDR
1567 analogue to the p-value¹⁵⁵.

1568
1569 In metabolomics, it is crucial to apply FDR correction methods to univariate results to ensure that
1570 the identified significant metabolites are not just statistical artifacts but reflect genuine biological
1571 differences. In all our univariate tests, we apply the BH metric to our p-values.

1572 1573 **Step 51: Install Packages for Univariate Analyses** ● **Timing 5 mins**

1574 Start by installing the packages necessary for this section: FSA⁸⁹ (v0.9.4), matrixStats⁹⁰ (v0.63.0).

1575 **3.4.1. Test for Normality**

1576 Testing for normality is often one of the first steps in univariate analysis and is crucial in deciding
1577 whether to use parametric or non-parametric tests. Parametric tests like t-test or ANOVA assume
1578 data follows a normal distribution, characterized by a symmetric bell-shaped curve with two key
1579 parameters: mean and standard deviation. Thus, before applying any statistical test, it is common
1580 to evaluate for normality with tests such as the Shapiro-Wilk test or the Kolmogorov-Smirnov test.
1581 Notably, Shapiro-Wilk is more suitable for small sample sizes ($N < 50$). Here, "normal" applies to
1582 the entire population, and not just the sample data. The resulting 'p value' from these tests only
1583 indicates the probability of the data to be sampled originating from a normal distribution. A
1584 graphical representation of testing normality of features is shown in **Figure 16**. Normality
1585 becomes less critical with large samples due to the Central Limit Theorem. In such cases,
1586 parametric tests can still be applied regardless of the normality. When the data does not follow a
1587 normal distribution, one can follow non-parametric tests, such as the Mann-Whitney U test or the
1588 Kruskal-Wallis test¹⁵⁶. In addition to this, to know more on normality assumptions, refer to **Box**
1589 **10**.

1590 In our pipeline, we conduct a normality test using two approaches: visual representations such as
1591 histograms and quantile-quantile plots (Q-Q plots), and the Shapiro-Wilk statistical test.

1592 1593 **Step 52: Normality Testing for One Feature**

1594 To illustrate how to test for normality, pick one feature and generate a Q-Q plot using the
1595 `qqnorm()` and `qqline()` functions. Then, perform a Shapiro-Wilk test using the
1596 `shapiro.test()` function. Additionally, demonstrate how log-transforming the data can
1597 improve normality.

1598
1599
1600

1601 **Step 53: Normality Testing for All Features**

1602 Perform a Shapiro-Wilk test for each feature and record the resulting p-values. Correct these p-
1603 values for false discovery rate (FDR) using the Benjamini & Hochberg method. If the adjusted p-
1604 value ('p_adj') is less than 0.05, reject the null hypothesis and consider the data to be non-
1605 normal. Tally up the features that fall under normal and non-normal distributions. If the majority of
1606 features are non-normal, consider using non-parametric tests for further analysis.

1607

1608 **Box 10 - Normality assumptions**

Besides normality, it is essential to consider two other critical assumptions when deciding between parametric and non-parametric tests: homogeneity of variances (homoscedasticity) and independence. Homoscedasticity demands that within-group variances are equal. If unequal (heteroscedasticity), it increases the chance of falsely identifying a “significant” result. Homoscedasticity can be evaluated graphically via boxplots or statistically via Levene’s and Bartlett’s tests. Here, the null hypothesis (H0) for these tests states that the within-group variances are equal. If the p-value is less than 0.05, it indicates a difference in population variances. The final assumption, ‘independence’, stipulates that the occurrence of one event does not influence the probability of another. In a metabolomic context, this implies that knowledge of one sample value does not predict another’s. However, these assumptions, particularly normality, are seldom fully met in real-world metabolomics datasets^{157,158}.

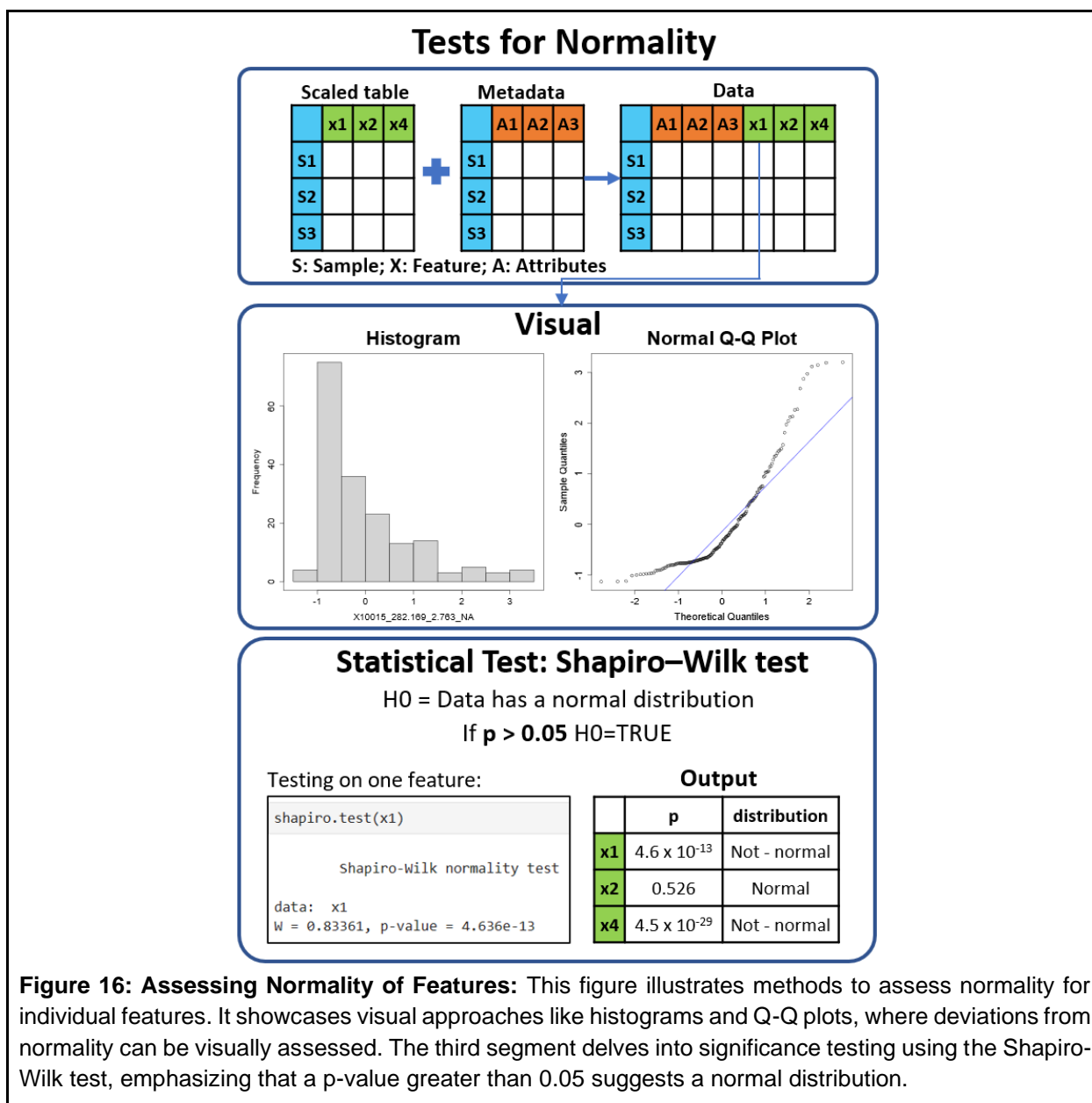


Figure 16: Assessing Normality of Features: This figure illustrates methods to assess normality for individual features. It showcases visual approaches like histograms and Q-Q plots, where deviations from normality can be visually assessed. The third segment delves into significance testing using the Shapiro-Wilk test, emphasizing that a p-value greater than 0.05 suggests a normal distribution.

1609

1610 3.4.2. Parametric tests

1611 3.4.2.1. ANOVA test

1612 The analysis of variance (ANOVA) is the statistical procedure used to test if there exists a
 1613 significant difference in the means of a dependent variable between three or more groups. As
 1614 opposed to a pair-wise comparison where we compare the means in a variable (i.e., $\mu_1=\mu_2$), in
 1615 the ANOVA we compare the means of several groups¹⁵⁹. For a deeper understanding of ANOVA,
 1616 please refer to **Box 11**. Furthermore, **Figure 17** offers a visual explanation of the ANOVA

1617 algorithm, detailing both the R code and a resulting plot that contrasts the F-statistic with p-values,
1618 highlighting significant features.

1619 **Step 54: Running ANOVA on one feature**

1620 *(User Input Required)*

1621 Here the user is prompted to enter the index number of the attribute for performing ANOVA. In
1622 the tutorial, we use 'ATTRIBUTE_Sample_Area'. The resulting ANOVA statistics are shown in a
1623 table format.

1624 **Step 55: Running ANOVA on all features**

- 1625 • For each metabolite feature, execute an ANOVA test within a for loop. The output for each
1626 feature is stored in a dataframe named ``anova_out``. The 'for loop' passes each feature
1627 column as the first argument of the `aov()` function against the selected attribute from the
1628 previous step ('ATTRIBUTE_Sample_Area'). This is because we are examining how a
1629 particular feature varies across different sample areas.
- 1630 • Tidy up the ANOVA output for each feature into a table using the `tidy()` function from
1631 the broom package.
- 1632 • Out of the two rows in the ANOVA summary table, select only the first row of this table
1633 (which contains the means, F-statistic, and p-value) and leave the second row consisting
1634 of the residuals.
- 1635 • Consolidate these rows into a single dataframe which contains the features, their
1636 corresponding p-values, their BH-corrected p-values, and their significance status in
1637 several columns. Features with a BH-corrected p-value (``anova_out$p_BH``) less than
1638 0.05 are considered significant.

1639 **Step 56: Subsetting Significant Features**

1640 Filter out the significant features for further examination. Display the count of significant and non-
1641 significant features.

1642 **Step 57: Visualize ANOVA Results**

1643 Sort the ``anova_out`` results by p-value and visualize the significant features using `ggplot()`.
1644 This involves plotting log-transformed F-Statistic values on the x-axis against negative logarithm
1645 of ``p_BH`` values on the y-axis. As F-Statistic and p-values can vary greatly, their log values offer
1646 easier visualization. To prevent clutter, limit the display to the names of the top 6 significant
1647 features.

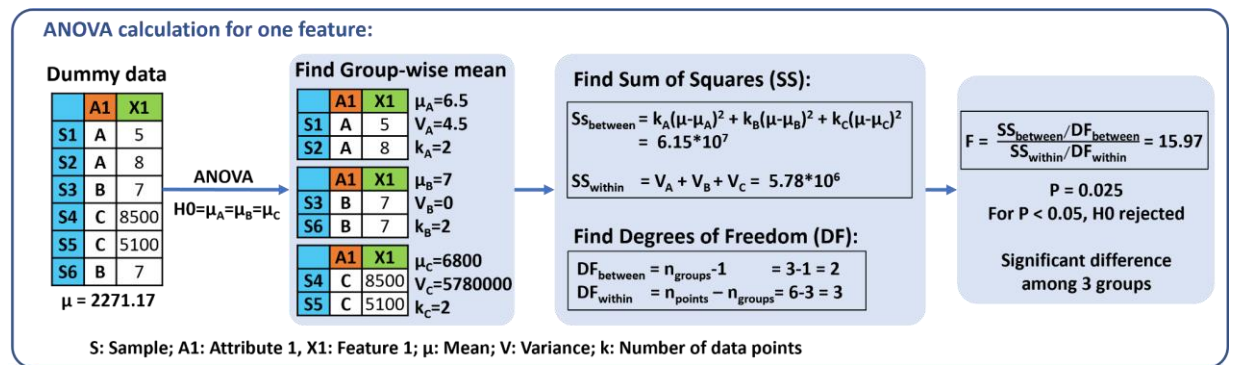
1648 **Step 58: Visualize Top Significant Metabolites**

1649 Generate boxplots for the top 4 significant metabolites to observe how their intensity levels differ
1650 across sampling sites. Extract these metabolites' data from the ``uni_data`` dataframe, which

1651 contains both feature intensities and metadata, and plot their intensities based on the sampling
 1652 sites. In our example, the higher intensities of these features in the 'Mission Bay' sample area
 1653 primarily account for the observed differences between sampling sites.

1654 **Box 11 - ANOVA**

If a pairwise test is used (e.g., a t-test), an increased probability of getting a false positive difference (Type I error) would be observed just by chance due to the effects of multiple comparisons¹⁶⁰. Instead, in the ANOVA test we can perform a single test to see if the observed differences are due to randomness or due to the grouping of the samples (e.g., origin, location, type of soil, etc.). The F-statistic is calculated using the sum of squares and the degrees of freedom (see **Figure 17**) and compared to a standard F-distribution to determine whether the differences among group means are greater than would be expected by chance. Importantly, the alternative hypothesis (i.e., where a difference exists between the means) is unpecific. This means that the test does not tell us where the difference(s) lie (e.g., if the difference is $\mu_A \neq \mu_B$ or $\mu_B \neq \mu_C$), it only tells us whether there exists a difference among all the means. The first assumption of the ANOVA test is the normality of population distribution and the homogeneity in their variances^{157,161}. Non-parametric tests should be used if these assumptions do not hold in the data of interest.

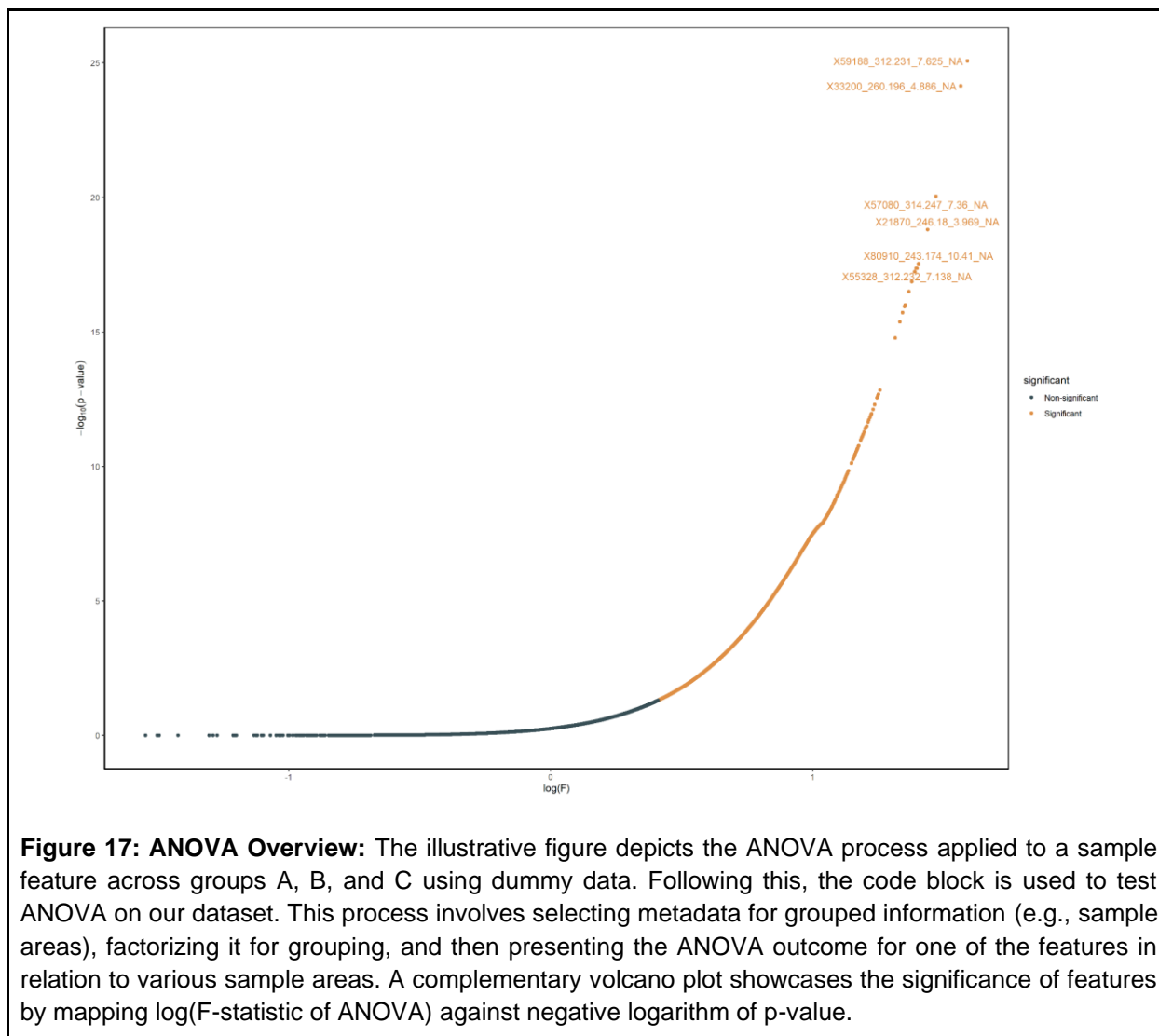


```
anova_group <- uni_metadata[,interested_attribute_anova]
anova_group <- as.factor(anova_group) # convert the attribute to 'factor' type

broom::tidy(aov(uni_data[,1] ~ anova_group)) #tidy summarizes the anova ouptut in a tibble
```

A tibble: 2 × 6

term	df	sumsq	meansq	statistic	p.value
<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
anova_group	6	7.560749	1.2601248	1.271597	0.2728111
Residuals	173	171.439251	0.9909783	NA	NA



1655

1656 3.4.2.2. Tukey's Honestly Significant Difference (HSD) test

1657 If the ANOVA test provides evidence that a difference indeed exists between the means of the
 1658 groups, the next step is to find between which groups the difference or differences exist. To do
 1659 this, we can conduct a Tukey HSD post hoc test used to compare multiple means in a single
 1660 analysis¹⁵⁷. Refer to **Box 12** for more information on Tukey's test. Additionally, **Figure 18** provides
 1661 a visual guide for applying the Tukey test, its implementation in R, and a resulting volcano plot
 1662 that highlights significant features from our pairwise comparison.

1663 Step 59: Perform Tukey HSD for a Significant Feature

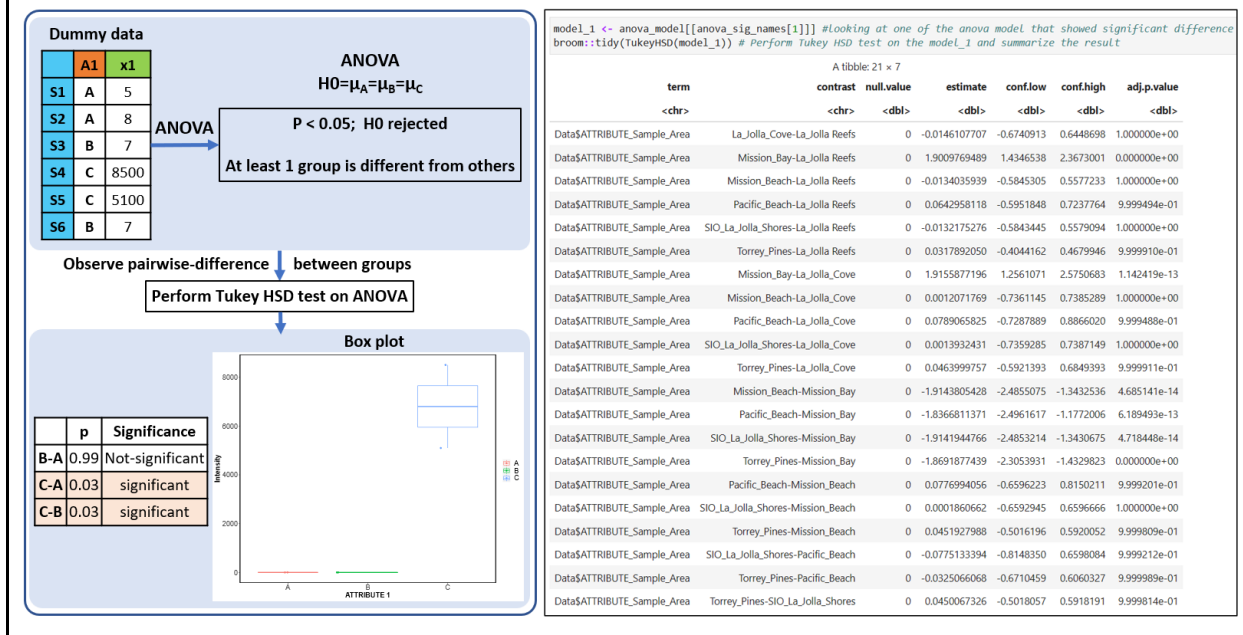
1664 First, we select a feature identified as significant in the ANOVA result, using
 1665 ``anova_sig_names`` generated in step 56. From the ANOVA output, we subset the data for

1666 this significant feature and conduct a Tukey HSD test. The output is a comprehensive table
 1667 providing an assessment of every possible pairwise group difference as shown in the figure.
 1668 To conduct a Tukey HSD test for all features, consider specifying just a one-pair comparison to
 1669 maintain simplicity. For instance, based on the ANOVA results, the sampling site 'Mission Bay'
 1670 appeared to significantly differ from others for the top four metabolites, hence we can focus on
 1671 the results from comparisons between 'Mission Bay' and another specific sampling site in the
 1672 subsequent step.

1673
 1674

Box 12 - Tukey's post hoc test

One of the goals of this test is to overcome the Type I error rate inflation of doing multiple comparisons¹⁵⁷. The most used post hoc test for ANOVA is Tukey's Honestly Significant Difference (HSD). To calculate the HSD between two means, a statistical distribution defined by Student (called the q distribution) is used which takes into account the number of means being compared¹⁶².



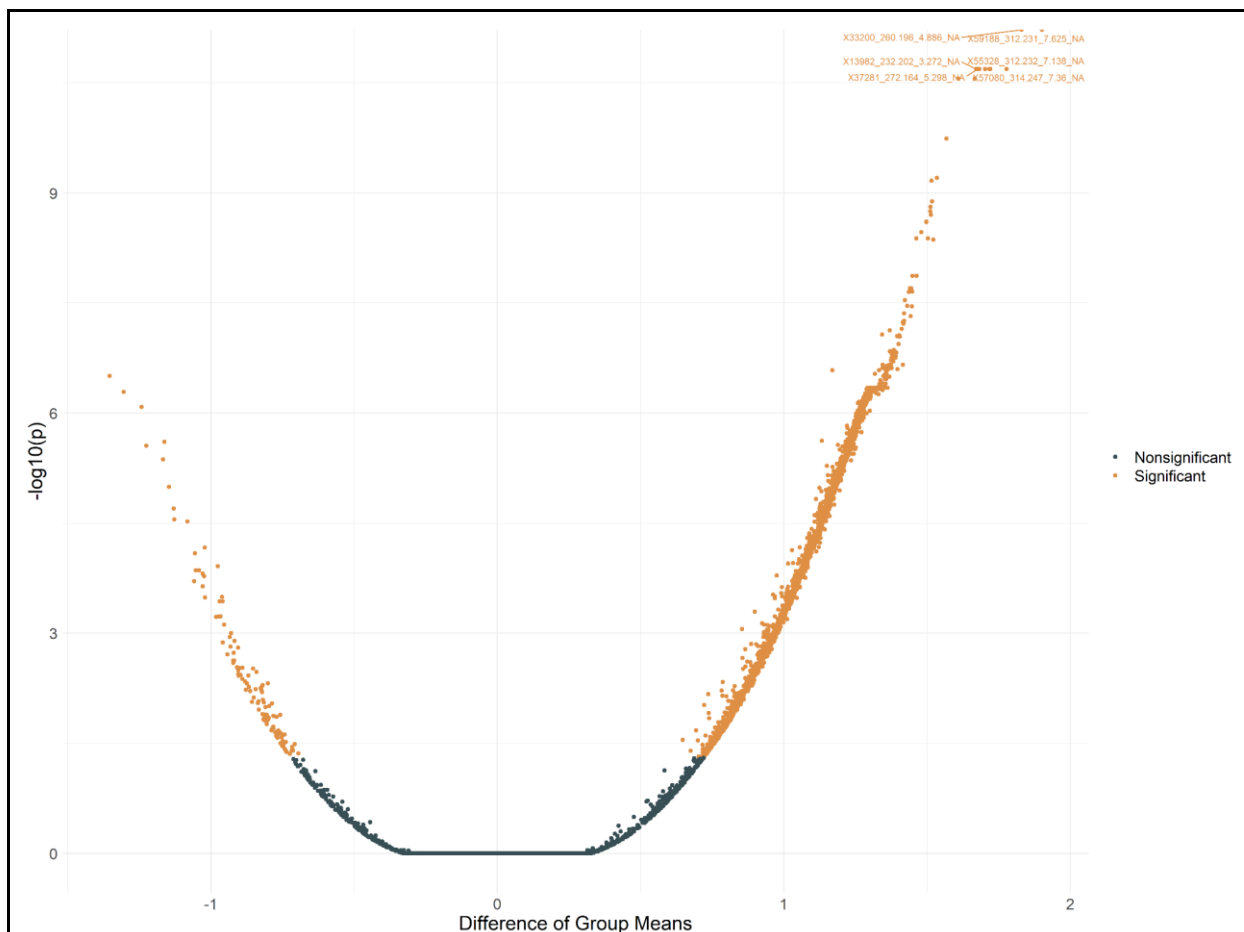


Figure 18: Overview of Tukey’s HSD Test: This figure starts with an illustration that showcases how ANOVA’s significance suggests that at least one group differs significantly from others, which then necessitates a further analysis through pairwise comparisons using Tukey’s HSD test. Alongside, we present the code block demonstrating the Tukey test applied to the first significant feature identified via ANOVA. Given the presence of 7 sample areas, the output presents p-values for all potential 21 pairwise comparisons. Having executed this for all ANOVA-significant features, we particularly highlighted comparisons between ‘Mission Bay’ and ‘La Jolla Reef’. The resulting significance is visualized via a volcano plot, where right-tailed features exhibit higher prevalence in ‘Mission Bay’, while left-tailed features dominate in ‘La Jolla Reef’.

1675 **Step 60: Perform Tukey HSD for All Significant Features**

1676 *(User Input Required)*

1677 Carry out a Tukey HSD test for all the significant features identified in the ANOVA. Then, filter the
 1678 results for the specific comparison such as ‘Mission Bay vs. La Jolla Reefs’. Here, users are
 1679 prompted to input the index number corresponding to their desired comparison from the ‘contrast’
 1680 column displayed in the previous step’s output. As a result of the Tukey test of this pairwise
 1681 interaction, p-values are produced for each feature. After applying the BH correction method,
 1682 features with corrected p-values (`output_tukey$p_BH < 0.05`) are highlighted as significantly
 1683 different between the selected sites.

1684 **Step 61: Count Significant Features**

1685 Determine how many features exhibit a significant difference between the chosen sites and how
1686 many do not.

1687 **Step 62: Visualize Results with a Volcano Plot**

1688 Create a volcano plot with `'-log(p_BH)'` on the y-axis and the group difference (`'estimate'`) on
1689 the x-axis. Display the names of the top findings on the plot to highlight the most significant
1690 differences between the chosen sites. Additionally, visualize the top 2 significant metabolites as
1691 boxplots from both extremes of the volcano plot (right and left tips) to clearly represent if the
1692 significant metabolite is upregulated or downregulated among the chosen sites.

1693 **3.4.2.3. T-tests**

1694 **Step 63: Select Attribute for T-Test Analysis**

1695 *(User Input Required)*

1696 A t-test is suitable for comparisons involving just two groups. Therefore, users should specify the
1697 attribute for the two distinct groups by providing the corresponding index number. For our
1698 example, we explore the metabolome's response to rainfall. Hence, we introduce an
1699 `'ATTRIBUTE_rainfall'` column, designating '1' for 'Jan-2018' (a high rainfall period) and '0' for
1700 the remaining months.

1701 **▲ CRITICAL:** This column addition caters to our dataset's context. Users with pre-existing binary
1702 attributes can skip this addition, while others may adjust this step to align with their data.

1703 **Step 64: Perform T-Test**

1704 Following the same steps as ANOVA (from steps 54 to 57), the `t.test()` function is used in
1705 place of `aov()` in this case. The final output is a dataframe `'ttest_output'` containing the
1706 significance of each feature for the two conditions under investigation.

1707 **Step 65: Plot T-Test Results**

1708 Visualize the t-test results using a volcano plot, with the `'estimate'` (difference in means of the
1709 two conditions for each feature) on the x-axis and `'-log(p_BH)'` on the y-axis. Additionally,
1710 visualize the top 2 significant metabolites as boxplots from both extremes of the volcano plot (right
1711 and left tips) to clearly represent if the significant metabolite is upregulated or downregulated for
1712 the chosen attribute.

1713 **▲ CRITICAL:** Unlike ANOVA, post-hoc tests are not needed for t-tests as there are only two
1714 conditions to compare. In ANOVA, when a feature is found to be significant, post-hoc tests help
1715 determine which specific groups show significant differences.

1716 For all the above tests, the respective significance values can be saved as a CSV table, and the
1717 plots can be saved in SVG, PDF, or PNG formats for further analysis or presentation.

1718 **3.4.3. Non-Parametric Tests**

1719 **3.4.3.1. Kruskal-Wallis Test**

1720 The Kruskal-Wallis test is a non-parametric statistical test used to compare three or more
1721 independent groups. It can be used when the assumptions of normality and equal variances are
1722 not met for performing an ANOVA¹⁶³. For more information on Kruskal-Wallis Test, refer to **Box**
1723 **13**. Additionally, **Figure 19** shows a visual explanation of the Kruskal-Wallis algorithm,
1724 accompanied by the R-code used to test a feature across various groups and determine its
1725 significance.

1726 **Step 66: Perform Kruskal-Wallis Test on one feature**

1727 *(User Input Required)*

1728 Begin by specifying the attribute for the Kruskal-Wallis (KW) test by entering its index number. In
1729 this tutorial, we opt for 'ATTRIBUTE_Sample_Area'. Then, apply the KW test on a single feature
1730 (the first feature in the `uni_data` dataframe) across different sample areas using the
1731 `kruskal.test()` function. Note that the `uni_data` dataframe originates from the
1732 `cleaned_data`, which we chose as the `Imp_s` scaled table (see Step 24). Summarize the
1733 output into a one-row table using the `tidy()` function from the broom package as shown in the
1734 figure.

1735 The steps for the Kruskal-Wallis test (steps 66 to 68) are structured similarly to the ANOVA steps
1736 (steps 54 to 57).

1737 **Step 67: Run Kruskal-Wallis Test for All Features**

- 1738
- 1739 • Just like in ANOVA (step 55), perform the Kruskal-Wallis test for each metabolite across
1740 different sample areas. Then, tidy up the output for each feature into a table using the
`tidy()` function.
 - 1741 • Combine these rows into a single dataframe containing features, their corresponding p-
1742 values, their BH-corrected p-values, and their significance status. Features with a BH-
1743 corrected p-value (`kruskal1_out$p_BH < 0.05`) less than 0.05 are considered
1744 significant.

1745 **Step 68: Filter Significant Features**

1746 Display the count of significant and non-significant features. Filter out the names of significant
1747 features for further analysis.

1748 **Step 69: Visualize Kruskal-Wallis Results**

1749 Similar to ANOVA results, we first sort the ``kruskall_out`` dataframe results by p-value and
1750 visualize the significant features using `ggplot()`. This involves plotting log-transformed K-
1751 Statistic values on the x-axis against `'-log(p_BH)'` on the y-axis. To prevent clutter, limit the
1752 display to the names of the top 6 significant features.

1753 **Step 70: Visualize Top Significant Metabolites of Kruskal-Wallis Results**

1754 Generate boxplots for the top 4 significant metabolites to observe how their intensity levels differ
1755 across sampling sites. Extract these metabolites' data from the ``uni_data`` dataframe, which
1756 contains both feature intensities and metadata, and plot their intensities based on the sampling
1757 sites.

1758 **Step 71: Compare Results from ANOVA and Kruskal-Wallis**

1759 We also suggest comparing the significant outcomes from both ANOVA and Kruskal-Wallis tests.
1760 Features yielding high scores in both tests indicate that the null hypothesis is rejected by both
1761 ANOVA and Kruskal-Wallis. This suggests that these features show significant differences across
1762 groups (in our case, across different sample areas). This comparison can help prioritize the
1763 features for further analysis.

1764 **Box 13 - Kruskal-Wallis test**

Although the Kruskal-Wallis test does not assume normality, it is expected that samples are random and independent and that the observations in each group come from populations with the same shape of distribution¹⁶³. As an extension of the Mann–Whitney U test (which is used to compare only two groups), it compares the median ranks of the groups, which are calculated by combining the ranks of all the observations across all groups and then taking their average¹⁶⁴. With this information, the K statistic can be calculated and compared to the chi-square distribution to accept or reject the null hypothesis (**Figure 19**). If the null hypothesis is rejected, the alternative hypothesis states that at least one group has a different median from the others.

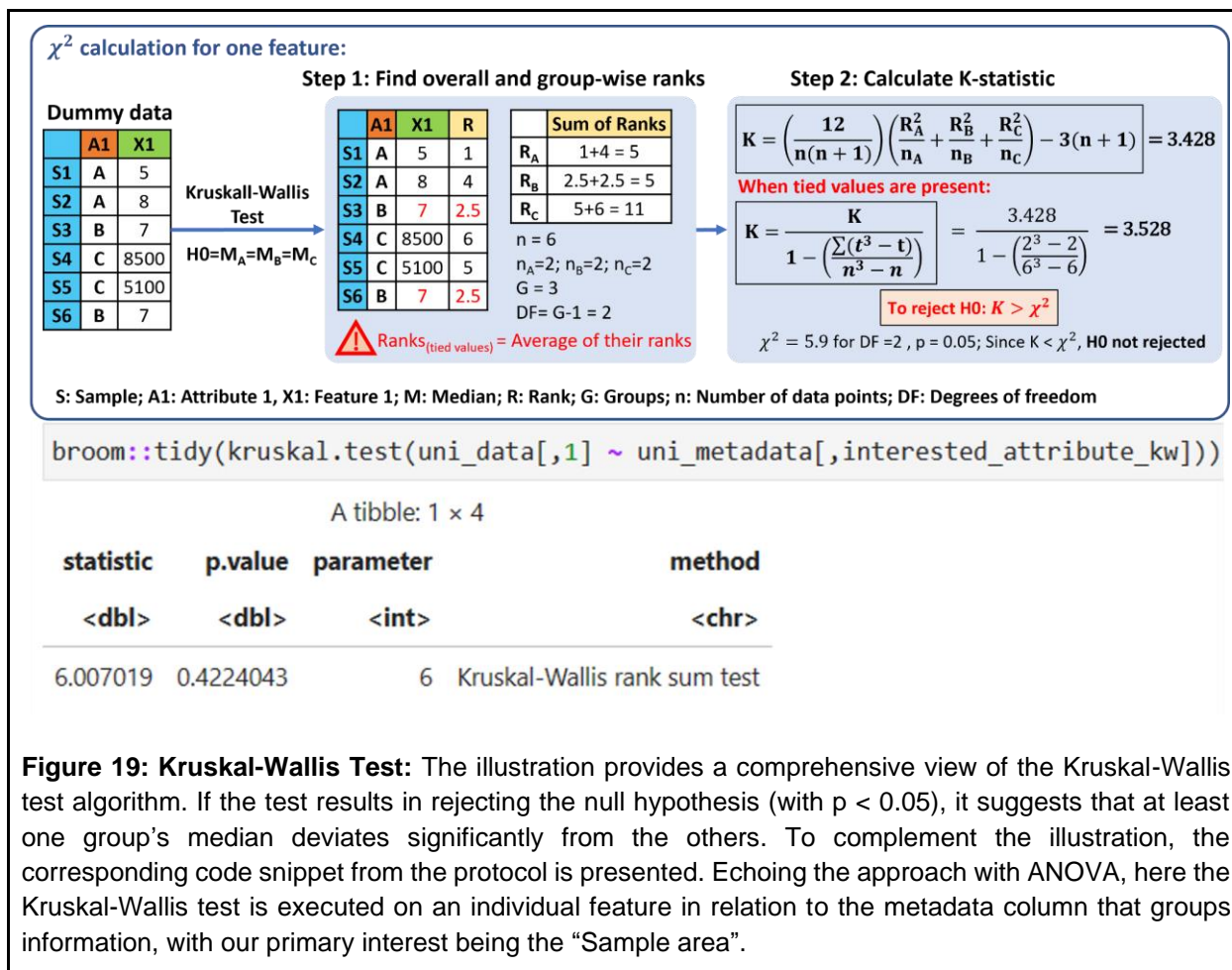


Figure 19: Kruskal-Wallis Test: The illustration provides a comprehensive view of the Kruskal-Wallis test algorithm. If the test results in rejecting the null hypothesis (with $p < 0.05$), it suggests that at least one group's median deviates significantly from the others. To complement the illustration, the corresponding code snippet from the protocol is presented. Echoing the approach with ANOVA, here the Kruskal-Wallis test is executed on an individual feature in relation to the metadata column that groups information, with our primary interest being the "Sample area".

1765

1766 3.4.3.2. Dunn's Post Hoc Test

1767 The Dunn statistical test is a non-parametric alternative to the Tukey HSD post hoc test to make
 1768 pairwise comparisons between multiple groups. The steps for Dunn's post hoc test (steps 72 to
 1769 75) are structured similarly to the Tukey HSD steps (steps 59 to 62). Refer to **Box 14** for more
 1770 information on Dunn's post hoc test. **Figure 20** shows a visual representation for applying the
 1771 Dunn test and its implementation in R.

1772 Step 72: Perform Dunn Test for a Significant Feature

1773 First, we select the first feature identified as significant in the KW test result, using
 1774 ``kw_sig_names`` generated in step 68. From the KW output, we subset the data for this
 1775 significant feature and conduct a Dunn test using `dunnTest()` function. The output is a
 1776 comprehensive table providing an assessment of every possible pairwise group difference as
 1777 shown in **Figure 20**.

1778 When conducting a Dunn test on all significant features, consider specifying just one pair
1779 interaction to maintain simplicity. Similar to the Tukey HSD test, here we will focus on the results
1780 from comparisons between 'Mission Bay' and 'La Jolla Reefs' in the subsequent step.

1781 **Step 73: Perform Dunn Test for All Significant Features**

1782 *(User Input Required)*

1783 Carry out a Dunn test for all the significant features identified in the Kruskal-Wallis test with BH
1784 correction for p-values. Then, filter the results for the specific interaction 'Mission Bay vs La Jolla
1785 Reefs'. To perform this, the user will be prompted to enter the index number corresponding to the
1786 desired comparison. This index number can be referenced from the table produced in the
1787 preceding step. Then, the Dunn Test result for those comparisons will be filtered for each feature
1788 showing the corrected p-values. The significance is assigned based on the corrected p-values
1789 (`dunn_output$P.adj < 0.05`) to identify the features that show a significant difference between
1790 these two sites.

1791 **Step 74: Count Significant Features**

1792 Determine how many features exhibit a significant difference between the chosen sites and how
1793 many do not.

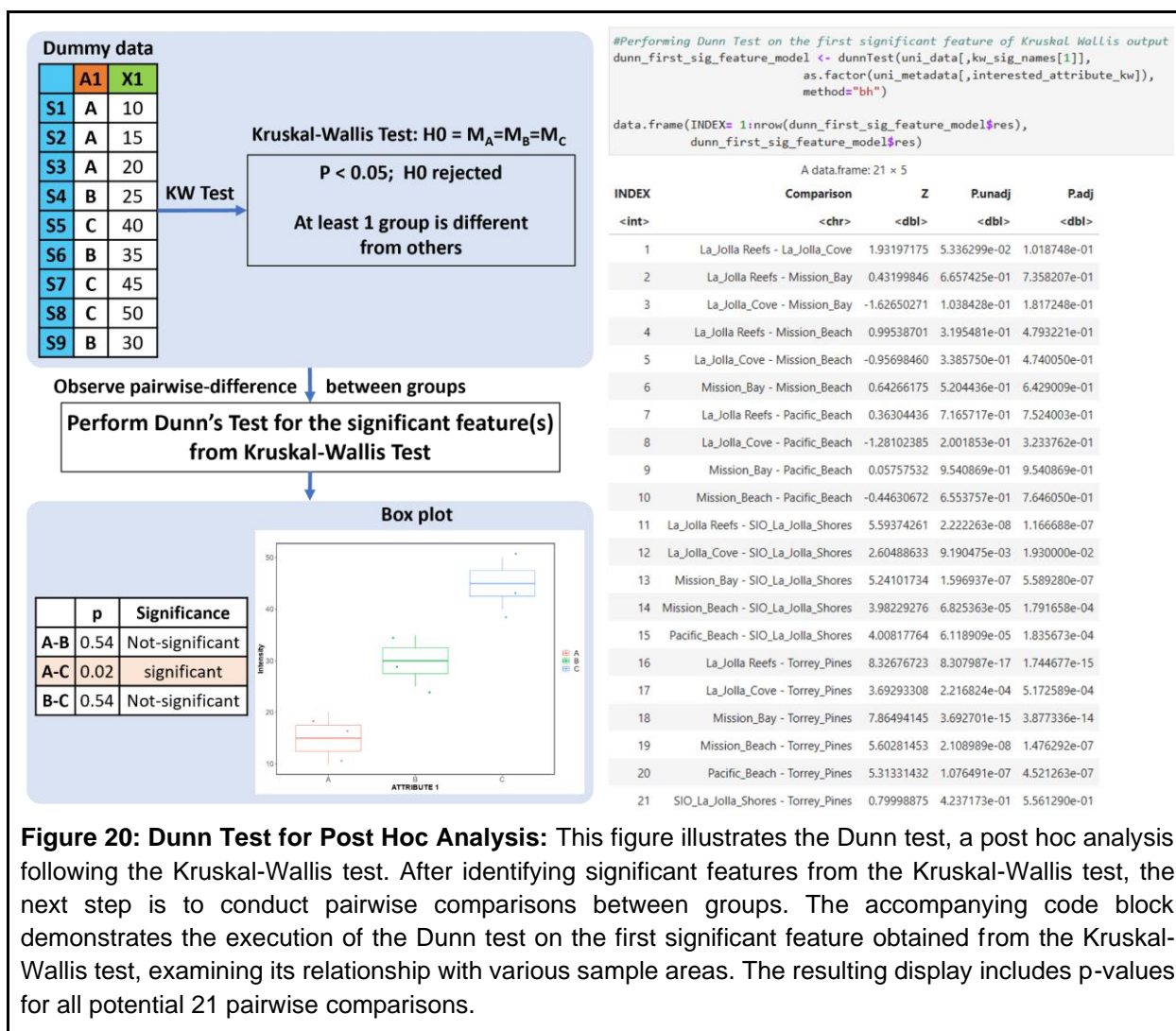
1794 **Step 75: Visualize Results with a Volcano Plot**

1795 Create a volcano plot with ' $-\log(p_{BH})$ ' on the y-axis and the Z statistic on the x-axis. Display
1796 the names of the top findings on the plot to highlight the most significant differences between the
1797 chosen sites. Additionally, visualize the top 2 significant metabolites as boxplots from both
1798 extremes of the volcano plot (right and left tips) to clearly represent if the significant metabolite is
1799 upregulated or downregulated for the chosen sites.

1800 For all the above tests, the respective significance values can be saved as a CSV table, and the
1801 plots can be saved in SVG, PDF, or PNG formats for further analysis or presentation.

1802 **Box 14 - Dunn test**

The Dunn statistical test is a non-parametric post-hoc test following Kruskal-Wallis test similar to the Tukey HSD post hoc test for ANOVA to make pairwise comparisons between multiple groups. Dunn's z-test approximation of the exact rank-sum test statistics is calculated with the mean rankings from the preceding Kruskal-Wallis test based on the differences in mean ranks for each group and, then, the p-value is calculated using a modified version of the BH correction to account for the type I error rate increase due to multiple comparisons¹⁶⁵ (**Figure 20**).



1803 4. Example Study

1804 Data Refinement and Annotation Insights

1805
 1806 In the example data, we investigated the coastal environments along the San Diego coastline
 1807 from Torrey Pines State Beach to Mission Bay, USA, during different dry and wet seasons. Refer
 1808 to **Figure 21A** for a spatial map of the sampling locations. The presumption was that post-rain
 1809 samples, influenced by runoff, would show increased pollutant levels. From FBMN analysis, we
 1810 identified 5521 LC-MS/MS features, which decreased to 4384 after removing blanks. The library
 1811 search against the GNPS spectral library via the FBMN workflow resulted in 92 annotated features
 1812 out of the 4384 features, and an additional analog search putatively annotated 104 features.
 1813 Expanding on this, we included additional data from October 2018, collected from the same sites
 1814 (no-rain period). for our pipeline evaluation. The dataset contained 180 samples from seven
 1815 different sites at three different time points (Dec 2017, Jan 2018, Oct 2018) and 2 PPL process
 1816 blanks for each sample time. From this extended dataset, we identified 11217 features, with 260

1817 GNPS library matches and 1991 analog matches. When focusing solely on December and
1818 January samples, the feature count surged to 10470, almost double the original count of 5521
1819 features, and 240 GNPS library hits and 1624 analog hits.

1820 To further expand our annotations, we used SIRIUS for in silico spectrum annotation on the
1821 extended dataset. We utilized the mgf file obtained from MZmine 3 and extended our SIRIUS
1822 analysis using tools like CANOPUS and CSI:FingerID. The SIRIUS result provided annotations
1823 for 8255 features, with annotations or compound names available for 5001 features. All 5001 of
1824 these features were further characterized by CSI:FingerID, which predicts molecular
1825 substructures and scores them based on the likelihood that the substructure belongs to the
1826 molecule. Leveraging the predictive capabilities of both SIRIUS and CSI:FingerID, we could infer
1827 the most probable molecular formulas. SIRIUS formula identifications were generated for 8885
1828 features, with 5411 of these having an explained intensity greater than 80%, marking them as
1829 reliable formulas. For compound class predictions, CANOPUS provided annotations for 8583
1830 features spanning various levels such as Kingdom, Superclass, Class, Subclass, and Level 5. On
1831 the other hand, the Natural Product Classifier (NPC) was used to determine if a compound is a
1832 natural product. These compound classes can be further explored in tools like Cytoscape for
1833 network visualization based on compound classes, or sub-setting of feature for subsequent
1834 statistics.

1835

1836 **Impact of Sequential Data Cleanup**

1837

1838 Contaminant features, especially those exceeding 30% peak area relative to the sample average,
1839 were flagged and removed, leaving us with 9,092 features. Our dataset showed 32% missing
1840 values out of 1,636,560 total entries, which were imputed between 1 and the lowest feature value
1841 (892). Petras *et al.* found significant organic matter chemotype shifts between December 2017
1842 and January 2018 samples, correlating with January's heavy rainfall¹⁴. Our extended dataset
1843 confirmed this, with a PCoA analysis revealing clear sample groupings by the sampling month as
1844 shown in **Figure 21B**. Post-blank removal intensified these groupings. Prior to data cleanup, no
1845 dispersion effect was apparent ($p > 0.05$), and PERMANOVA attributed 31% of the variance to
1846 sampling months. After removing blanks, however, a dispersion effect emerged. This dispersion
1847 effect and explained variance in PERMANOVA are likely due to the removal of background
1848 features, thus reflecting the true water sample chemotypes for each month. Upon examining the
1849 PCoA after imputation, individual clusters appeared closer together, though January samples
1850 exhibited some dispersion. This spread within January samples became more pronounced after
1851 normalization and scaling.

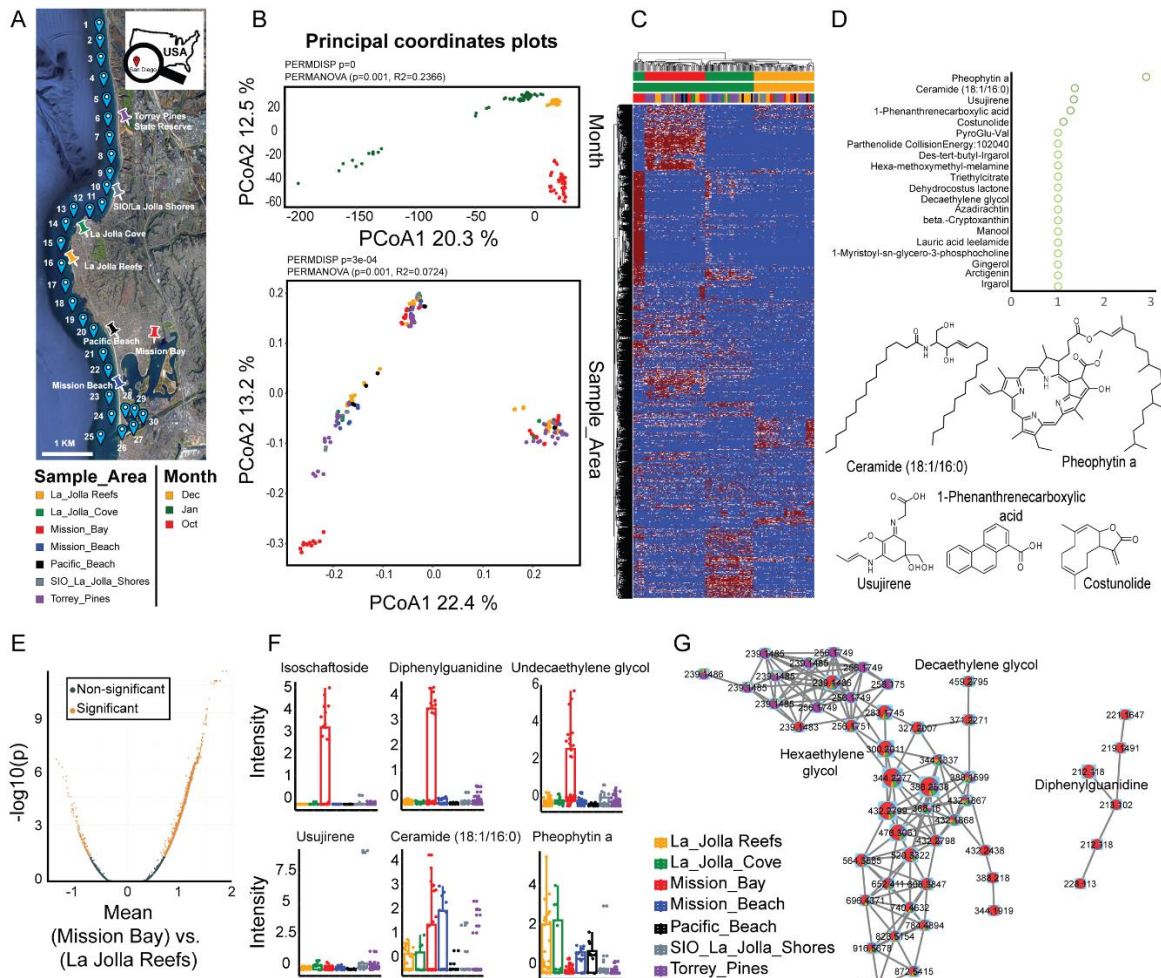
1852

1853 **Multivariate Analysis: Diving into Site-Specific Variations**

1854

1855 Using PERMANOVA on the scaled-imputed data, we identified a significant clustering by months,
1856 attributing 34% of variance to the sampling time ($P < 0.05$, Adonis $R^2 = 0.34$). Sample locations,
1857 however, explained only 7% of the variance. Upon deeper exploration at the metabolic profiles
1858 across these sampling locations, January's variance was more prominent in Mission Bay,
1859 especially post-rainfall, due to its nutrient-rich status, potentially from increased runoff through the
1860 San Diego River. This distinction is evident in the PCoA plot in **Figure 21B**. Our data showed

1861 Mission Bay's pre-rainfall samples were similar to other sites, but post-rain samples in January
 1862 diverged — a pattern absent in December 2017 and October 2018 samples. We could also
 1863 observe some clear patterns in the heatmap depicted in **Figure 21C**. Color transitions from blue
 1864 (0 intensity) to red (1 intensity) highlight feature intensity variations. Many features were found in
 1865 higher intensities in October samples compared to December and January samples. Mission Bay
 1866 samples from January (in red) and a subset from Torrey Pines (in blue) displayed increased
 1867 feature intensities. This aligns well with our initial hypothesis. Alongside this, we performed a
 1868 random forest classification considering sampling sites.
 1869



1870
 1871 **Figure 21: Anticipated results:** A) Spatial map pinpointing sampling sites; B) Principal coordinate plots
 1872 delineating differences by sampling month and location; C) Heatmap displaying scaled feature intensities;
 1873 D) Top 20 annotated drivers for temporal changes identified via Random Forest, with structures of the top
 1874 5 metabolites shown; E) Volcano plot illustrating the Tukey test comparison between Mission Bay and La
 1875 Jolla Reefs samples, with features deemed significant in the ANOVA Sample area-based test used for this
 1876 post-hoc analysis; F) Box plots illustrating feature intensities across various sampling locations. The top
 1877 row presents the foremost 3 annotated significant outcomes post-Tukey test, accompanied by their
 1878 molecular structures. Conversely, the second row highlights the top 3 significant outputs as identified by
 1879 Random Forest; G) Molecular Networks of significant features (diphenylguanidine and polyethylene glycols,
 1880 highlight related molecules with similar spatial patterns as indicated through the pie charts on top of the
 1881 network nodes.

1882 Random Forest Exploration: Prioritizing Key Drivers

1883 Utilizing a Random Forest model with 500 trees and 500 permutations, we attained a 68.3%
1884 prediction accuracy for the samples. By location, accuracy ranged from 87.5% (Torrey Pines) to
1885 16.7% (Pacific Beach). The confusion matrix in **Table 4** provides insights into these results,
1886 revealing that misclassifications were often between neighboring sites, likely due to the close 300-
1887 meter spacing between the sampling locations. Our model highlighted 438 significant features
1888 (based on 'Mean Decrease Accuracy p value'). Of these, seven matched GNPS libraries and 96
1889 were analog hits. Examining the violin plot results of RF, top features, like those with library IDs
1890 91372 and 90597 (both sharing the same analog name), were mainly concentrated in Mission
1891 Bay and La Jolla Reefs. These concentrations began low at Torrey Pines, peaked at Cove and
1892 Reef, and saw another spike in Mission Bay. Similar patterns emerged for features like theaflavin
1893 digallatae (ID 91133). Some features, such as IDs 33200 and 53617, were notably elevated in
1894 Mission Bay alone. Certain compounds from previously reported research, such as *m/z*
1895 1129.3145 (analog name: benzyl-tetradecyl-dimethylammonium) specific to January samples,
1896 were also detected in our study, but their significance was marginal ($p = 0.08$) and was
1897 predominantly seen in Torrey Pines. Several compounds reported in the original study such as
1898 irgarol, recognized for their pollution potential and unique spatial patterns, were also explored in
1899 our dataset. **Figure 21D** visualizes the top 20 annotated drivers for site-specific changes as
1900 identified via Random Forest, highlighting the structures of the top 5 metabolites. In summary, our
1901 extended data set enhances the Random Forest analysis, offering a detailed understanding of
1902 chemotype differences across coastal areas and reaffirming the conclusions of the original study.
1903

1904 **Table 4: Confusion Matrix of Random Forest Classification**

1905 The confusion matrix shows how many samples from each group were correctly predicted. Taking
1906 the first row as an example: out of 36 samples from La Jolla Reefs, 25 were accurately identified.
1907 The remaining samples were misclassified as follows: 1 as Mission Bay, 1 as Mission Beach, 5
1908 as Pacific Beach, and 4 as SIO La Jolla Shores. The column labeled 'pct.correct' represents the
1909 percentage of samples that were correctly classified for a given group. The columns 'LCI 0.95'
1910 and 'UCI 0.95' denote the lower and upper bounds of the 95% confidence interval for each group,
1911 respectively. The 'overall' row at the bottom indicates the model's total prediction accuracy, which
1912 stands at 68.3% for this dataset.
1913

	La Jolla Reefs	La Jolla Cove	Mission Bay	Mission Beach	Pacific Beach	SIO La Jolla Shores	Torrey Pines	pct.correct	LCI 0.95	UCI 0.95
La Jolla Reefs	25	0	1	1	5	4	0	69.4	51.89	83.7
La Jolla Cove	0	10	0	0	0	2	0	83.3	51.59	97.9
Mission Bay	4	0	23	7	2	0	0	63.9	46.22	79.2
Mission Beach	0	0	0	15	3	0	0	83.3	58.58	96.4
Pacific Beach	6	0	1	3	2	0	0	16.7	2.09	48.4
SIO La Jolla Shores	2	0	0	1	0	6	9	33.3	13.34	59
Torrey Pines	0	0	0	0	0	6	42	87.5	74.75	95.3
Overall	NA	NA	NA	NA	NA	NA	NA	68.3	61	75.1

1914

1915

1916 **Univariate Analysis Insights**

1917
1918 In our univariate analysis of 9092 features, both ANOVA and the Kruskal-Wallis test were utilized.
1919 However, the Kruskal-Wallis test was considered more apt due to the non-parametric nature of
1920 our dataset. The Kruskal-Wallis test highlighted 1258 significant features, including irgarol, an
1921 antifouling agent used on boats. Conversely, ANOVA pinpointed 1554 significant features, with
1922 many features having a pronounced abundance in Mission Bay compared to other sites. Notably,
1923 one of the features corresponded to hexaethylene glycol from the NIST14 database and several
1924 features matched to nanoethylene glycols, which fall under the PEGs (Polyethylene glycols)
1925 category. Another notable find was an analog match to sporidesmolide 2, previously identified in
1926 the base study.

1927 Building on the ANOVA results, Tukey's HSD test was used to highlight pairwise differences.
1928 Given the pronounced abundance of many features in Mission Bay, we compared it with La Jolla
1929 Reefs for further insights. The significant and non-significant features from this test are visualized
1930 in the volcano plot in **Figure 21E**. Notably, compounds like 1,2-diphenylguanidine (used in metal
1931 detection and rubber vulcanization) and nanoethylene glycol were significantly higher in Mission
1932 Bay. In contrast, La Jolla Reefs had a higher presence of the natural product 'pheophytin a' at
1933 various retention times (RT 11.312, 11.022). The top row in **Figure 21F** displays the intensities of
1934 the top three annotated results from the Tukey test across the sampling locations using box plots,
1935 each paired with its corresponding molecular structure. Interestingly, 'pheophytin a' was also more
1936 abundant in La Jolla Reefs in subsequent Dunn Tests post-Kruskal-Wallis. Furthermore, irgarol
1937 was consistently found to be more abundant in Mission Bay in both tests.

1938 These findings align with and reinforce the initial observations, validating the robustness of our
1939 analytical workflow.

1940

1941 **Integration of Molecular Networking Results**

1942
1943 After the statistical analysis of the FBMN results and prioritization of features that drive the
1944 chemical differences between the sampling sites, we further investigate related compounds,
1945 through the molecular networks. **Figure 21G** shows the networks of diphenylguanidine and
1946 polyethylene glycols, indicating that many of the structurally related features of those compounds
1947 show similar spatial distribution, with the highest abundance in Mission Bay, as indicated through
1948 the pie charts on top of the network nodes. These results show nicely how the statistical
1949 prioritization and further structure-based (in our case, based on MS/MS similarity) can work hand
1950 in hand to structure the observed chemical space. Besides investigating the networks after the
1951 statistical interrogation, one can also make use of the scores obtained from the different tests and
1952 visualize those in the network. For example, the fold change and p-values from the univariate
1953 analysis or mean decreased accuracies from the supervised multivariate analysis can be imported
1954 as new attribute to the networks with tools such as Cytoscape to combine visual and statistical
1955 prioritization directly in the network.

1956 5. Conclusion

1957 In this protocol, we provide a comprehensive data clean-up and statistics pipeline for the analysis
1958 of non-targeted metabolomics data. Our protocol spans from initial data conversion, blank
1959 removal, imputation, and normalization/scaling to uni- and multivariate statistics and data
1960 interpretation. While our outlined workflow is as detailed and structured as possible, which should
1961 provide a comprehensive analysis solution for many biological questions, it is important to point
1962 out that there is not a universal solution that fits every scenario. We emphasize the importance of
1963 transparency in reporting details on every step of the metabolomics pipeline, such as providing
1964 the specific normalization methods, explaining the distance metrics in multivariate analysis, or
1965 specifying parameters like the number of trees in a Random Forest model. Furthermore, in
1966 relation to our case study, the sharing of feature detection and annotation settings and batch files
1967 further augments reproducibility. Together, with open data deposition, the above steps ensure
1968 both transparency and reproducibility of metabolomics experiments.

1969 We would also like to stress again that cataloging and identifying statistically significant
1970 metabolites is just the beginning. To fully understand the relationships between metabolites and
1971 the underlying biological processes, additional experiments and orthogonal verification are
1972 typically required. Once the statistical results are studied, techniques such as pathway enrichment
1973 analyses can illuminate the multifaceted relationships between metabolites and the biological
1974 processes they are entwined with. When specific compounds are of particular interest, targeted
1975 metabolomics stands as a powerful next step.

1976 In summary, we anticipate that our Hitchhiker's Guide to statistical analysis of FBMN results will
1977 provide both a theoretical and practical resource for scientists working with non-targeted
1978 metabolomics data. For novices in the field, the scripts, app and detailed step-to-step protocol
1979 provide a starting point with a set of statistical analysis solutions for many biological questions,
1980 whereas experts may accelerate parts of their statistical workflows.

1981 6. References

- 1982 1. Vailati-Riboni, M., Palombo, V. & Loor, J. J. What Are Omics Sciences? in *Periparturient*
1983 *Diseases of Dairy Cows: A Systems Biology Approach* (ed. Ametaj, B. N.) 1–7 (Springer
1984 International Publishing, 2017). doi:10.1007/978-3-319-43033-1_1.
- 1985 2. Patti, G. J., Yanes, O. & Siuzdak, G. Metabolomics: the apogee of the omics trilogy. *Nat. Rev.*
1986 *Mol. Cell Biol.* **13**, 263–269 (2012).
- 1987 3. Dayalan, S., Xia, J., Spicer, R. A., Salek, R. & Roessner, U. Metabolome Analysis. in
1988 *Encyclopedia of Bioinformatics and Computational Biology* (eds. Ranganathan, S., Gribskov,
1989 M., Nakai, K. & Schönbach, C.) 396–409 (Academic Press, 2019). doi:10.1016/B978-0-12-
1990 809633-8.20251-3.
- 1991 4. Tolstikov, V., Moser, A. J., Sarangarajan, R., Narain, N. R. & Kiebish, M. A. Current Status of
1992 Metabolomic Biomarker Discovery: Impact of Study Design and Demographic Characteristics.
1993 *Metabolites* **10**, 224 (2020).
- 1994 5. de Jonge, N. F. *et al.* Good practices and recommendations for using and benchmarking
1995 computational metabolomics metabolite annotation tools. *Metabolomics* **18**, 103 (2022).

- 1996 6. Nothias, L.-F. *et al.* Feature-based molecular networking in the GNPS analysis environment.
1997 *Nat. Methods* **17**, 905–908 (2020).
- 1998 7. Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global Natural
1999 Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).
- 2000 8. Ottosson, F. *et al.* Effects of Long-Term Storage on the Biobanked Neonatal Dried Blood Spot
2001 Metabolome. *J. Am. Soc. Mass Spectrom.* **34**, 685–694 (2023).
- 2002 9. Dantas Machado, A. C. *et al.* Portosystemic shunt placement reveals blood signatures for the
2003 development of hepatic encephalopathy through mass spectrometry. *Nat. Commun.* **14**, 5303
2004 (2023).
- 2005 10. Xie, H.-F. *et al.* Feature-based molecular networking analysis of the metabolites produced
2006 by in vitro solid-state fermentation reveals pathways for the bioconversion of epigallocatechin
2007 gallate. *J. Agric. Food Chem.* **68**, 7995–8007 (2020).
- 2008 11. Berlanga-Clavero, M. V. *et al.* *Bacillus subtilis* biofilm matrix components target seed oil
2009 bodies to promote growth and anti-fungal resistance in melon. *Nat. Microbiol.* **7**, 1001–1015
2010 (2022).
- 2011 12. Raheem, D. J., Tawfike, A. F., Abdelmohsen, U. R., Edrada-Ebel, R. & Fitzsimmons-
2012 Thoss, V. Application of metabolomics and molecular networking in investigating the chemical
2013 profile and antitrypanosomal activity of British bluebells (*Hyacinthoides non-scripta*). *Sci. Rep.*
2014 **9**, 2547 (2019).
- 2015 13. Pendergraft, M. A. *et al.* Bacterial and Chemical Evidence of Coastal Water Pollution from
2016 the Tijuana River in Sea Spray Aerosol. *Environ. Sci. Technol.* **57**, 4071–4081 (2023).
- 2017 14. Petras, D. *et al.* Non-targeted tandem mass spectrometry enables the visualization of
2018 organic matter chemotype shifts in coastal seawater. *Chemosphere* **271**, 129450 (2021).
- 2019 15. Stincone, P. *et al.* Evaluation of Data-Dependent MS/MS Acquisition Parameters for Non-
2020 Targeted Metabolomics and Molecular Networking of Environmental Samples: Focus on the Q
2021 Exactive Platform. *Anal. Chem.* **95**, 12673–12682 (2023).
- 2022 16. Wegley Kelly, L. *et al.* Distinguishing the molecular diversity, nutrient content, and
2023 energetic potential of exometabolomes produced by macroalgae and reef-building corals. *Proc.*
2024 *Natl. Acad. Sci.* **119**, e2110283119 (2022).
- 2025 17. Mannocho-Russo, H. *et al.* Microbiomes and metabolomes of dominant coral reef primary
2026 producers illustrate a potential role for immunolipids in marine symbioses. *Commun. Biol.* **6**,
2027 896 (2023).
- 2028 18. Shaffer, J. P. *et al.* Standardized multi-omics of Earth's microbiomes reveals microbial and
2029 metabolite diversity. *Nat. Microbiol.* **7**, 2128–2150 (2022).
- 2030 19. Molina-Santiago, C. *et al.* Chemical interplay and complementary adaptative strategies
2031 toggle bacterial antagonism and co-existence. *Cell Rep.* **36**, (2021).
- 2032 20. Reher, R. *et al.* Native metabolomics identifies the rivulariapeptolide family of protease
2033 inhibitors. *Nat. Commun.* **13**, 4619 (2022).
- 2034 21. Aron, A. T. *et al.* Native mass spectrometry-based metabolomics identifies metal-binding
2035 compounds. *Nat. Chem.* **14**, 100–109 (2022).
- 2036 22. Behnsen, J. *et al.* Siderophore-mediated zinc acquisition enhances enterobacterial
2037 colonization of the inflamed gut. *Nat. Commun.* **12**, 7016 (2021).
- 2038 23. Pang, Z. *et al.* MetaboAnalyst 5.0: narrowing the gap between raw spectra and functional
2039 insights. *Nucleic Acids Res.* **49**, W388–W396 (2021).

- 2040 24. Pang, Z. *et al.* Using MetaboAnalyst 5.0 for LC–HRMS spectra processing, multi-omics
2041 integration and covariate adjustment of global metabolomics data. *Nat. Protoc.* **17**, 1735–1761
2042 (2022).
- 2043 25. Cajka, T. & Fiehn, O. Toward Merging Untargeted and Targeted Methods in Mass
2044 Spectrometry-Based Metabolomics and Lipidomics. *Anal. Chem.* **88**, 524–545 (2016).
- 2045 26. Alder, L., Greulich, K., Kempe, G. & Vieth, B. Residue analysis of 500 high priority
2046 pesticides: Better by GC–MS or LC–MS/MS? *Mass Spectrom. Rev.* **25**, 838–865 (2006).
- 2047 27. Díaz-Cruz, M. S., López de Alda, M. J., López, R. & Barceló, D. Determination of
2048 estrogens and progestogens by mass spectrometric techniques (GC/MS, LC/MS and
2049 LC/MS/MS). *J. Mass Spectrom.* **38**, 917–923 (2003).
- 2050 28. Michely, J. A., Helfer, A. G., Brandt, S. D., Meyer, M. R. & Maurer, H. H. Metabolism of
2051 the new psychoactive substances N,N-diallyltryptamine (DALT) and 5-methoxy-DALT and their
2052 detectability in urine by GC–MS, LC–MSn, and LC–HR–MS–MS. *Anal. Bioanal. Chem.* **407**,
2053 7831–7842 (2015).
- 2054 29. Di Masi, S. *et al.* HPLC-MS/MS method applied to an untargeted metabolomics approach
2055 for the diagnosis of “olive quick decline syndrome”. *Anal. Bioanal. Chem.* **414**, 465–473 (2022).
- 2056 30. Reveglia, P. *et al.* Untargeted and Targeted LC-MS/MS Based Metabolomics Study on In
2057 Vitro Culture of Phaeoacremonium Species. *J. Fungi* **8**, 55 (2022).
- 2058 31. Baig, F., Pechlaner, R. & Mayr, M. Caveats of Untargeted Metabolomics for Biomarker
2059 Discovery*. *J. Am. Coll. Cardiol.* **68**, 1294–1296 (2016).
- 2060 32. Xiao, J. F., Zhou, B. & Ressom, H. W. Metabolite identification and quantitation in LC-
2061 MS/MS-based metabolomics. *TrAC Trends Anal. Chem.* **32**, 1–14 (2012).
- 2062 33. Blaženović, I. *et al.* Comprehensive comparison of in silico MS/MS fragmentation tools of
2063 the CASMI contest: database boosting is needed to achieve 93% accuracy. *J.*
2064 *Cheminformatics* **9**, 32 (2017).
- 2065 34. Blaženović, I., Kind, T., Ji, J. & Fiehn, O. Software Tools and Approaches for Compound
2066 Identification of LC-MS/MS Data in Metabolomics. *Metabolites* **8**, 31 (2018).
- 2067 35. Dührkop, K., Shen, H., Meusel, M., Rousu, J. & Böcker, S. Searching molecular structure
2068 databases with tandem mass spectra using CSI:FingerID. *Proc. Natl. Acad. Sci.* **112**, 12580–
2069 12585 (2015).
- 2070 36. Böcker, S., Letzel, M. C., Lipták, Z. & Pervukhin, A. SIRIUS: decomposing isotope patterns
2071 for metabolite identification†. *Bioinformatics* **25**, 218–224 (2009).
- 2072 37. Stravs, M. A., Dührkop, K., Böcker, S. & Zamboni, N. MSNovelist: de novo structure
2073 generation from mass spectra. *Nat. Methods* **19**, 865–870 (2022).
- 2074 38. Aron, A. T. *et al.* Reproducible molecular networking of untargeted mass spectrometry
2075 data using GNPS. *Nat. Protoc.* **15**, 1954–1991 (2020).
- 2076 39. Schmid, R. *et al.* Ion identity molecular networking for mass spectrometry-based
2077 metabolomics in the GNPS environment. *Nat. Commun.* **12**, 3832 (2021).
- 2078 40. Dunn, W. B. *et al.* Procedures for large-scale metabolic profiling of serum and plasma
2079 using gas chromatography and liquid chromatography coupled to mass spectrometry. *Nat.*
2080 *Protoc.* **6**, 1060–1083 (2011).
- 2081 41. Silva, A. M., Cordeiro-da-Silva, A. & Coombs, G. H. Metabolic Variation during
2082 Development in Culture of *Leishmania donovani* Promastigotes. *PLoS Negl. Trop. Dis.* **5**,
2083 e1451 (2011).

- 2084 42. Martínez-Sena, T. *et al.* Monitoring of system conditioning after blank injections in
2085 untargeted UPLC-MS metabolomic analysis. *Sci. Rep.* **9**, 9822 (2019).
- 2086 43. Steffen Heuckeroth, Tito Damiani, Aleksandr Smirnov, Olena Mokshyna, Corinna Brungs,
2087 Ansgar Korf, Joshua David Smith, Paolo Stincone, Nicola Dreolin, Louis-Félix Nothias, Tuulia
2088 Hyötyläinen, Matej Orešič, Uwe Karst, Pieter C. Dorrestein, Daniel Petras, Xiuxia Du, Justin
2089 J.J. van der Hoof, Robin Schmid, Tomáš Pluskal. Mass spectrometry data processing in
2090 MZmine 3: feature detection and annotation. *Nat. Protoc.* **Under Review**, (2023).
- 2091 44. Raynie, D. The Vital Role of Blanks in Sample Preparation. *LCGC N. Am.* **36**, 494–497
2092 (2018).
- 2093 45. Kessner, D., Chambers, M., Burke, R., Agus, D. & Mallick, P. ProteoWizard: open source
2094 software for rapid proteomics tools development. *Bioinformatics* **24**, 2534–2536 (2008).
- 2095 46. Hulstaert, N. *et al.* ThermoRawFileParser: Modular, Scalable, and Cross-Platform RAW
2096 File Conversion. *J. Proteome Res.* **19**, 537–542 (2020).
- 2097 47. Smith, C. A., Want, E. J., O’Maille, G., Abagyan, R. & Siuzdak, G. XCMS: Processing
2098 Mass Spectrometry Data for Metabolite Profiling Using Nonlinear Peak Alignment, Matching,
2099 and Identification. *Anal. Chem.* **78**, 779–787 (2006).
- 2100 48. Kuhl, C., Tautenhahn, R., Böttcher, C., Larson, T. R. & Neumann, S. CAMERA: An
2101 Integrated Strategy for Compound Spectra Extraction and Annotation of Liquid
2102 Chromatography/Mass Spectrometry Data Sets. *Anal. Chem.* **84**, 283–289 (2012).
- 2103 49. Schmid, R. *et al.* Integrative analysis of multimodal mass spectrometry data in MZmine 3.
2104 *Nat. Biotechnol.* **41**, 447–449 (2023).
- 2105 50. Tsugawa, H. *et al.* A lipidome atlas in MS-DIAL 4. *Nat. Biotechnol.* **38**, 1159–1163 (2020).
- 2106 51. Pfeuffer, J. *et al.* OpenMS – A platform for reproducible analysis of mass spectrometry
2107 data. *J. Biotechnol.* **261**, 142–148 (2017).
- 2108 52. Gloaguen, Y., Kirwan, J. A. & Beule, D. Deep Learning-Assisted Peak Curation for Large-
2109 Scale LC-MS Metabolomics. *Anal. Chem.* **94**, 4930–4937 (2022).
- 2110 53. Chetnik, K., Petrick, L. & Pandey, G. MetaClean: a machine learning-based classifier for
2111 reduced false positive peak detection in untargeted LC–MS metabolomics data. *Metabolomics*
2112 **16**, 117 (2020).
- 2113 54. El Abiead, Y., Milford, M., Salek, R. M. & Koellensperger, G. mzRAPP: a tool for reliability
2114 assessment of data pre-processing in non-targeted metabolomics. *Bioinformatics* **37**, 3678–
2115 3680 (2021).
- 2116 55. Tautenhahn, R., Böttcher, C. & Neumann, S. Highly sensitive feature detection for high
2117 resolution LC/MS. *BMC Bioinformatics* **9**, 504 (2008).
- 2118 56. Dührkop, K. *et al.* SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite
2119 structure information. *Nat. Methods* **16**, 299–302 (2019).
- 2120 57. Sumner, L. W. *et al.* Proposed minimum reporting standards for chemical analysis.
2121 *Metabolomics* **3**, 211–221 (2007).
- 2122 58. Liu, L.-L. *et al.* Molecular networking-based for the target discovery of potent
2123 antiproliferative polycyclic macrolactam ansamycins from *Streptomyces cacaoi* subsp.
2124 *asoensis*. *Org. Chem. Front.* **7**, 4008–4018 (2020).
- 2125 59. Sedio, B. E., Boya P., C. A. & Rojas Echeverri, J. C. A protocol for high-throughput,
2126 untargeted forest community metabolomics using mass spectrometry molecular networks.
2127 *Appl. Plant Sci.* **6**, e1033 (2018).

- 2128 60. Quinn, R. A. *et al.* Molecular Networking As a Drug Discovery, Drug Metabolism, and
2129 Precision Medicine Strategy. *Trends Pharmacol. Sci.* **38**, 143–154 (2017).
- 2130 61. Pluskal, T., Castillo, S., Villar-Briones, A. & Orešič, M. MZmine 2: Modular framework for
2131 processing, visualizing, and analyzing mass spectrometry-based molecular profile data. *BMC*
2132 *Bioinformatics* **11**, 395 (2010).
- 2133 62. Djoumbou Feunang, Y. *et al.* ClassyFire: automated chemical classification with a
2134 comprehensive, computable taxonomy. *J. Cheminformatics* **8**, 61 (2016).
- 2135 63. Kim, H. W. *et al.* NPClassifier: A Deep Neural Network-Based Structural Classification
2136 Tool for Natural Products. *J. Nat. Prod.* **84**, 2795–2807 (2021).
- 2137 64. Dührkop, K. *et al.* Systematic classification of unknown metabolites using high-resolution
2138 fragmentation mass spectra. *Nat. Biotechnol.* **39**, 462–471 (2021).
- 2139 65. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data
2140 science using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 2141 66. Davidson, R. L., Weber, R. J. M., Liu, H., Sharma-Oates, A. & Viant, M. R. Galaxy-M: a
2142 Galaxy workflow for processing and analyzing direct infusion and liquid chromatography mass
2143 spectrometry-based metabolomics data. *GigaScience* **5**, 10 (2016).
- 2144 67. Giacomoni, F. *et al.* Workflow4Metabolomics: a collaborative research infrastructure for
2145 computational metabolomics. *Bioinformatics* **31**, 1493–1495 (2015).
- 2146 68. Kontou, E. E. *et al.* UmetaFlow: an untargeted metabolomics workflow for high-throughput
2147 data processing and analysis. *J. Cheminformatics* **15**, 52 (2023).
- 2148 69. Rohart, F., Gautier, B., Singh, A. & Lê Cao, K.-A. mixOmics: An R package for 'omics
2149 feature selection and multiple data integration. *PLoS Comput. Biol.* **13**, e1005752 (2017).
- 2150 70. Chong, J. & Xia, J. MetaboAnalystR: an R package for flexible and reproducible analysis
2151 of metabolomics data. *Bioinformatics* **34**, 4313–4314 (2018).
- 2152 71. Tiffany, C. R. & Bäuml, A. J. omu, a Metabolomics count data analysis tool for intuitive
2153 figures and convenient metadata collection. *Microbiol. Resour. Announc.* **8**, 10.1128/mra.
2154 00129-19 (2019).
- 2155 72. Han, X. & Liang, L. metabolomicsR: a streamlined workflow to analyze metabolomic data
2156 in R. *Bioinforma. Adv.* **2**, vbac067 (2022).
- 2157 73. Fernández-Albert, F., Llorach, R., Andrés-Lacueva, C. & Perera, A. An R package to
2158 analyse LC/MS metabolomic data: MAIT (Metabolite Automatic Identification Toolkit).
2159 *Bioinformatics* **30**, 1937–1939 (2014).
- 2160 74. Kohler, D. *et al.* MSstats Version 4.0: Statistical Analyses of Quantitative Mass
2161 Spectrometry-Based Proteomic Experiments with Chromatography-Based Quantification at
2162 Scale. *J. Proteome Res.* **22**, 1466–1482 (2023).
- 2163 75. Riquelme, G., Zabalegui, N., Marchi, P., Jones, C. M. & Monge, M. E. A Python-Based
2164 Pipeline for Preprocessing LC–MS Data for Untargeted Metabolomics Workflows. *Metabolites*
2165 **10**, 416 (2020).
- 2166 76. Di Guida, R. *et al.* Non-targeted UHPLC-MS metabolomic data processing methods: a
2167 comparative investigation of normalisation, missing value imputation, transformation and
2168 scaling. *Metabolomics* **12**, 93 (2016).
- 2169 77. Watrous, J. *et al.* Mass spectral molecular networking of living microbial colonies. *Proc.*
2170 *Natl. Acad. Sci.* **109**, E1743–E1752 (2012).

- 2171 78. Hoffmann, M. A. *et al.* Assigning confidence to structural annotations from mass spectra
2172 with COSMIC. 2021.03.18.435634 Preprint at <https://doi.org/10.1101/2021.03.18.435634>
2173 (2021).
- 2174 79. Rinker, T. *et al.* pacman: Package Management Tool. (2019).
- 2175 80. Wickham, H. *et al.* Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).
- 2176 81. Kluyver, T., Angerer, P. & Schulz, J. IRdisplay: ‘Jupyter’ Display Machinery. (2022).
- 2177 82. Cacciatore, S., Luchinat, C. & Tenori, L. Knowledge discovery by accuracy maximization.
2178 *Proc. Natl. Acad. Sci.* **111**, 5117–5122 (2014).
- 2179 83. Kassambara, A. & Mundt, F. Extract and Visualize the Results of Multivariate Data
2180 Analyses [R package factoextra version 1.0.7]. in (2020).
- 2181 84. Oksanen, J. *et al.* vegan: Community Ecology Package. (2022).
- 2182 85. Gu, Z. Complex heatmap visualization. *iMeta* **1**, e43 (2022).
- 2183 86. Galili, T. dendextend: an R package for visualizing, adjusting and comparing trees of
2184 hierarchical clustering. *Bioinforma. Oxf. Engl.* **31**, 3718–3720 (2015).
- 2185 87. Charrad, M., Ghazzali, N., Boiteau, V. & Niknafs, A. NbClust: An R Package for
2186 Determining the Relevant Number of Clusters in a Data Set. *J. Stat. Softw.* **61**, 1–36 (2014).
- 2187 88. Archer, E. rfPermute: Estimate Permutation p-Values for Random Forest Importance
2188 Metrics. (2023).
- 2189 89. Ogle, D. H., Doll, J. C., Wheeler, A. P. & dunnTest(), A. D. (Provided base functionality
2190 of. FSA: Simple Fisheries Stock Assessment Methods. (2023).
- 2191 90. Bengtsson, H. *et al.* matrixStats: Functions that Apply to Rows and Columns of Matrices
2192 (and to Vectors). (2023).
- 2193 91. Xiao [aut, N., cre, Cook, J., Jégousse, C. & Li, M. ggsci: Scientific Journal and Sci-Fi
2194 Themed Color Palettes for ‘ggplot2’. (2023).
- 2195 92. Wilke, C. O. cowplot: Streamlined Plot Theme and Plot Annotations for ‘ggplot2’. (2020).
- 2196 93. Wickham, H. *et al.* svglite: An ‘SVG’ Graphics Device. (2023).
- 2197 94. Liu, Q. *et al.* Addressing the batch effect issue for LC/MS metabolomics data in data
2198 preprocessing. *Sci. Rep.* **10**, 13856 (2020).
- 2199 95. Yue, Y., Bao, X., Jiang, J. & Li, J. Evaluation and correction of injection order effects in
2200 LC-MS/MS based targeted metabolomics. *J. Chromatogr. B* **1212**, 123513 (2022).
- 2201 96. Livera, A. M. D. *et al.* Statistical Methods for Handling Unwanted Variation in
2202 Metabolomics Data. *Anal. Chem.* **87**, 3606–3615 (2015).
- 2203 97. Reese, S. E. *et al.* A new statistic for identifying batch effects in high-throughput genomic
2204 data that uses guided principal component analysis. *Bioinformatics* **29**, 2877–2883 (2013).
- 2205 98. Burton, L. *et al.* Instrumental and experimental effects in LC–MS-based metabolomics. *J.*
2206 *Chromatogr. B* **871**, 227–235 (2008).
- 2207 99. Gregori, J. *et al.* Batch effects correction improves the sensitivity of significance tests in
2208 spectral counting-based comparative discovery proteomics. *J. Proteomics* **75**, 3938–3951
2209 (2012).
- 2210 100. Thonusin, C. *et al.* Evaluation of intensity drift correction strategies using MetaboDrift, a
2211 normalization tool for multi-batch metabolomics data. *J. Chromatogr. A* **1523**, 265–274 (2017).
- 2212 101. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression
2213 data using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).

- 2214 102. Deng, K. *et al.* WavelCA: A novel algorithm to remove batch effects for large-scale
2215 untargeted metabolomics data based on wavelet analysis. *Anal. Chim. Acta* **1061**, 60–69
2216 (2019).
- 2217 103. Wehrens, R. *et al.* Improved batch correction in untargeted MS-based metabolomics.
2218 *Metabolomics* **12**, 88 (2016).
- 2219 104. Kuligowski, J., Sánchez-Illana, Á., Sanjuán-Herráez, D., Vento, M. & Quintás, G. Intra-
2220 batch effect correction in liquid chromatography-mass spectrometry using quality control
2221 samples and support vector regression (QC-SVRC). *The Analyst* **140**, 7810–7817 (2015).
- 2222 105. Luan, H., Ji, F., Chen, Y. & Cai, Z. statTarget: A streamlined tool for signal drift correction
2223 and interpretations of quantitative mass spectrometry-based omics data. *Anal. Chim. Acta*
2224 **1036**, 66–72 (2018).
- 2225 106. Rong, Z. *et al.* NormAE: Deep Adversarial Learning Model to Remove Batch Effects in
2226 Liquid Chromatography Mass Spectrometry-Based Metabolomics Data. *Anal. Chem.* **92**,
2227 5082–5090 (2020).
- 2228 107. Dmitrenko, A., Reid, M. & Zamboni, N. Regularized adversarial learning for normalization
2229 of multi-batch untargeted metabolomics data. *Bioinformatics* **39**, btad096 (2023).
- 2230 108. Tokareva, A. O. *et al.* Normalization methods for reducing interbatch effect without quality
2231 control samples in liquid chromatography-mass spectrometry-based studies. *Anal. Bioanal.*
2232 *Chem.* **413**, 3479–3486 (2021).
- 2233 109. Cleary, J. L., Luu, G. T., Pierce, E. C., Dutton, R. J. & Sanchez, L. M. BLANKA: an
2234 Algorithm for Blank Subtraction in Mass Spectrometry of Complex Biological Samples. *J. Am.*
2235 *Soc. Mass Spectrom.* **30**, 1426–1434 (2019).
- 2236 110. Lawson, T. N. *et al.* msPurity: Automated Evaluation of Precursor Ion Purity for Mass
2237 Spectrometry-Based Fragmentation in Metabolomics. *Anal. Chem.* **89**, 2432–2439 (2017).
- 2238 111. Schiffman, C. *et al.* Data-adaptive pipeline for filtering and normalizing metabolomics data.
2239 387365 Preprint at <https://doi.org/10.1101/387365> (2018).
- 2240 112. Carobene, A., Braga, F., Roraas, T., Sandberg, S. & Bartlett, W. A. A systematic review
2241 of data on biological variation for alanine aminotransferase, aspartate aminotransferase and γ -
2242 glutamyl transferase. *Clin. Chem. Lab. Med. CCLM* **51**, 1997–2007 (2013).
- 2243 113. Schiffman, C. *et al.* Filtering procedures for untargeted LC-MS metabolomics data. *BMC*
2244 *Bioinformatics* **20**, 334 (2019).
- 2245 114. Broadhurst, D. *et al.* Guidelines and considerations for the use of system suitability and
2246 quality control samples in mass spectrometry assays applied in untargeted clinical
2247 metabolomic studies. *Metabolomics* **14**, 72 (2018).
- 2248 115. Wei, R. *et al.* Missing Value Imputation Approach for Mass Spectrometry-based
2249 Metabolomics Data. *Sci. Rep.* **8**, 663 (2018).
- 2250 116. Do, K. T. *et al.* Characterization of missing values in untargeted MS-based metabolomics
2251 data and evaluation of missing data handling strategies. *Metabolomics* **14**, 128 (2018).
- 2252 117. Gorrochategui, E., Jaumot, J., Lacorte, S. & Tauler, R. Data analysis strategies for
2253 targeted and untargeted LC-MS metabolomic studies: Overview and workflow. *TrAC Trends*
2254 *Anal. Chem.* **82**, 425–442 (2016).
- 2255 118. Wulff, J. E. & Mitchell, M. W. A Comparison of Various Normalization Methods for LC/MS
2256 Metabolomics Data. *Adv. Biosci. Biotechnol.* **9**, 339–351 (2018).

- 2257 119. Dieterle, F., Ross, A., Schlotterbeck, G. & Senn, H. Probabilistic Quotient Normalization
2258 as Robust Method to Account for Dilution of Complex Biological Mixtures. Application in 1H
2259 NMR Metabonomics. *Anal. Chem.* **78**, 4281–4290 (2006).
- 2260 120. Li, B. *et al.* Performance Evaluation and Online Realization of Data-driven Normalization
2261 Methods Used in LC/MS based Untargeted Metabolomics Analysis. *Sci. Rep.* **6**, 38881 (2016).
- 2262 121. Scholz, M., Gatzek, S., Sterling, A., Fiehn, O. & Selbig, J. Metabolite fingerprinting:
2263 detecting biological features by independent component analysis. *Bioinformatics* **20**, 2447–
2264 2454 (2004).
- 2265 122. Qannari, E. M., Wakeling, I., Courcoux, P. & MacFie, H. J. H. Defining the underlying
2266 sensory dimensions. *Food Qual. Prefer.* **11**, 151–154 (2000).
- 2267 123. van den Berg, R. A., Hoefsloot, H. C., Westerhuis, J. A., Smilde, A. K. & van der Werf, M.
2268 J. Centering, scaling, and transformations: improving the biological information content of
2269 metabolomics data. *BMC Genomics* **7**, 142 (2006).
- 2270 124. Khalheim, O. M. Scaling of analytical data. *Anal. Chim. Acta* **177**, 71–79 (1985).
- 2271 125. Kasprzak, E. M. & Lewis, K. E. Pareto analysis in multiobjective optimization using the
2272 collinearity theorem and scaling method. *Struct. Multidiscip. Optim.* **22**, 208–218 (2001).
- 2273 126. Keenan, M. R. & Kotula, P. G. Accounting for Poisson noise in the multivariate analysis of
2274 ToF-SIMS spectrum images. *Surf. Interface Anal.* **36**, 203–212 (2004).
- 2275 127. Gromski, P. S. *et al.* A tutorial review: Metabolomics and partial least squares-discriminant
2276 analysis – a marriage of convenience or a shotgun wedding. *Anal. Chim. Acta* **879**, 10–23
2277 (2015).
- 2278 128. Mendez, K. M., Reinke, S. N. & Broadhurst, D. I. A comparative evaluation of the
2279 generalised predictive ability of eight machine learning algorithms across ten clinical
2280 metabolomics data sets for binary classification. *Metabolomics* **15**, 150 (2019).
- 2281 129. GOWER, J. C. Some distance properties of latent root and vector methods used in
2282 multivariate analysis. *Biometrika* **53**, 325–338 (1966).
- 2283 130. Xu, Y. *et al.* Application of Dissimilarity Indices, Principal Coordinates Analysis, and Rank
2284 Tests to Peak Tables in Metabolomics of the Gas Chromatography/Mass Spectrometry of
2285 Human Sweat. *Anal. Chem.* **79**, 5633–5641 (2007).
- 2286 131. Nguyen, L. H. & Holmes, S. Ten quick tips for effective dimensionality reduction. *PLOS*
2287 *Comput. Biol.* **15**, e1006907 (2019).
- 2288 132. Jäggi, C., Wirth, T. & Baur, B. Genetic variability in subpopulations of the asp viper (*Vipera*
2289 *aspis*) in the Swiss Jura mountains: implications for a conservation strategy. *Biol. Conserv.* **94**,
2290 69–77 (2000).
- 2291 133. Pinheiro, H. P., de Souza Pinheiro, A. & Sen, P. K. Comparison of genomic sequences
2292 using the Hamming distance. *J. Stat. Plan. Inference* **130**, 325–339 (2005).
- 2293 134. Lozupone, C. & Knight, R. UniFrac: a New Phylogenetic Method for Comparing Microbial
2294 Communities. *Appl. Environ. Microbiol.* **71**, 8228–8235 (2005).
- 2295 135. Brejnrod, A. *et al.* Implementations of the chemical structural and compositional similarity
2296 metric in R and Python. 546150 Preprint at <https://doi.org/10.1101/546150> (2019).
- 2297 136. Tripathi, A. *et al.* Chemically informed analyses of metabolomics mass spectrometry data
2298 with Qemistree. *Nat. Chem. Biol.* **17**, 146–151 (2021).
- 2299 137. Ramette, A. Multivariate analyses in microbial ecology. *FEMS Microbiol. Ecol.* **62**, 142–
2300 160 (2007).

- 2301 138. Koenig, J. E. *et al.* Succession of microbial consortia in the developing infant gut
2302 microbiome. *Proc. Natl. Acad. Sci.* **108**, 4578–4585 (2011).
- 2303 139. Paliy, O. & Shankar, V. Application of multivariate statistical techniques in microbial
2304 ecology. *Mol. Ecol.* **25**, 1032–1057 (2016).
- 2305 140. Efron, B. Bootstrap Methods: Another Look at the Jackknife. in *Breakthroughs in Statistics:
2306 Methodology and Distribution* (eds. Kotz, S. & Johnson, N. L.) 569–593 (Springer, 1992).
2307 doi:10.1007/978-1-4612-4380-9_41.
- 2308 141. Desu, M. M. & Raghavarao, D. *Nonparametric Statistical Methods For Complete and
2309 Censored Data*. (CRC Press, 2003).
- 2310 142. Anderson, M. J. A new method for non-parametric multivariate analysis of variance.
2311 *Austral Ecol.* **26**, 32–46 (2001).
- 2312 143. Anderson, M. J. & Walsh, D. C. I. PERMANOVA, ANOSIM, and the Mantel test in the face
2313 of heterogeneous dispersions: What null hypothesis are you testing? *Ecol. Monogr.* **83**, 557–
2314 574 (2013).
- 2315 144. Tibshirani, R., Walther, G. & Hastie, T. Estimating the number of clusters in a data set via
2316 the gap statistic. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **63**, 411–423 (2001).
- 2317 145. Wilkinson, L. & Friendly, M. The History of the Cluster Heat Map. *Am. Stat.* **63**, 179–184
2318 (2009).
- 2319 146. Wu, W. & Noble, W. S. Genomic data visualization on the Web. *Bioinformatics* **20**, 1804–
2320 1805 (2004).
- 2321 147. Benton, P. H. *et al.* An Interactive Cluster Heat Map to Visualize and Explore
2322 Multidimensional Metabolomic Data. *Metabolomics Off. J. Metabolomic Soc.* **11**, 1029–1034
2323 (2015).
- 2324 148. Ren, S., Hinzman, A. A., Kang, E. L., Szczesniak, R. D. & Lu, L. J. Computational and
2325 statistical analysis of metabolomics data. *Metabolomics* **11**, 1492–1513 (2015).
- 2326 149. Liebal, U. W., Phan, A. N. T., Sudhakar, M., Raman, K. & Blank, L. M. Machine Learning
2327 Applications for Mass Spectrometry-Based Metabolomics. *Metabolites* **10**, 243 (2020).
- 2328 150. Breiman, L. (out-of-bag estimates). (1996).
- 2329 151. Breiman, L. Random Forests. *Mach. Learn.* **45**, 5–32 (2001).
- 2330 152. Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T. & Zeileis, A. Conditional variable
2331 importance for random forests. *BMC Bioinformatics* **9**, 307 (2008).
- 2332 153. Archer, K. J. & Kimes, R. V. Empirical characterization of random forest variable
2333 importance measures. *Comput. Stat. Data Anal.* **52**, 2249–2260 (2008).
- 2334 154. Jafari, M. & Ansari-Pour, N. Why, When and How to Adjust Your P Values? *Cell J. Yakhteh*
2335 **20**, 604–607 (2019).
- 2336 155. Korthauer, K. *et al.* A practical guide to methods controlling false discoveries in
2337 computational biology. *Genome Biol.* **20**, 118 (2019).
- 2338 156. Mishra, P. *et al.* Descriptive Statistics and Normality Tests for Statistical Data. *Ann. Card.*
2339 *Anaesth.* **22**, 67–72 (2019).
- 2340 157. Vinaixa, M. *et al.* A Guideline to Univariate Statistical Analysis for LC/MS-Based
2341 Untargeted Metabolomics-Derived Data. *Metabolites* **2**, 775–795 (2012).
- 2342 158. Riffenburgh, R. H. & Gillen, D. L. *Statistics in Medicine*. (Academic Press, 2020).
- 2343 159. Xia, Y. & Sun, J. Hypothesis Testing and Statistical Analysis of Microbiome. *Genes Dis.*
2344 **4**, 138–148 (2017).

- 2345 160. Sato, T. Type I and Type II Error in Multiple Comparisons. *J. Psychol.* **130**, 293–302
2346 (1996).
- 2347 161. Bathke, A. The ANOVA F test can still be used in some balanced designs with unequal
2348 variances and nonnormal data. *J. Stat. Plan. Inference* **126**, 413–422 (2004).
- 2349 162. Abdi, H. & Williams, L. Newman-Keuls Test and Tukey Test. *Encycl. Res. Des.* (2010).
- 2350 163. Ostertagová, E., Ostertag, O. & Kováč, J. Methodology and Application of the Kruskal-
2351 Wallis Test. *Appl. Mech. Mater.* **611**, 115–120 (2014).
- 2352 164. Hecke, T. V. Power study of anova versus Kruskal-Wallis test. *J. Stat. Manag. Syst.* **15**,
2353 241–247 (2012).
- 2354 165. Dinno, A. Nonparametric Pairwise Multiple Comparisons in Independent Groups using
2355 Dunn’s Test. *Stata J. Promot. Commun. Stat. Stata* **15**, 292–300 (2015).
- 2356

2357 7. Data and Code Sharing

2358 The FBMN results are available under the following URL:

2359 <https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=b661d12ba88745639664988329c1363e>

2360 Raw and processed that is available through the MassIVE repository (MSV000090156) and
2361 Zenodo (<https://zenodo.org/records/10051610>). All code and software is available through GitHub
2362 under the following link <https://github.com/Functional-Metabolomics-Lab/FBMN-STATS>.

2363 8. Acknowledgment

2364 We thank Libera Lo Presti for critical reading of the manuscript. We thank Greg Caporaso for
2365 guidance on preparing the QIIME2 plugins. DP, CM, and HPL were supported by the Deutsche
2366 Forschungsgemeinschaft (DFG) through the CMFI Cluster of Excellence (EXC 2124) and DP and
2367 CM, were supported by the DFG through the Collaborative Research Center CellMap (TRR 261).
2368 KD was supported by the DFG (BO 1910/23). PS was supported by the European Union’s Horizon
2369 Europe research and innovation programme through a Marie Skłodowska-Curie fellowship No.
2370 101108450 MeStaLeM. TP was supported by the Czech Science Foundation (GA CR) grant 21-
2371 11563M and by the European Union’s Horizon 2020 research and innovation programme under
2372 Marie Skłodowska-Curie grant agreement No. 891397. TD was supported by the MSCA
2373 Fellowships CZ (OP JAK) grant CZ.02.01.01/00/22_010/0002733. MW was supported by the
2374 National Institutes of Health (NIH) with grants 1U24DK133658-01, NIH 1R03DE032437-01, and
2375 UC Riverside startup funding. EEK was supported by grants of the Novo Nordisk Foundation
2376 [NNF20CC0035580, NNF16OC0021746]. YW was supported by NIH 1R03DE032437-01. CB
2377 was supported by the Czech Academy of Sciences (CAS PPLZ) L200552251. FO was supported
2378 by FAPESP 2022/14603-8. JB was supported by Deutsches Zentrum für Infektionsforschung
2379 (DZIF). EEK was supported by grants of the Novo Nordisk Foundation (NNF20CC0035580,
2380 NNF16OC0021746).

2381 9. Author Contribution

2382 AKPS, FO, FR, ME, and DP conceptualized the protocol. YE, SZ, JS, RS advised on the concept
2383 and statistical test. AKPS, AW, FO, FR, MN, JB, EEK, JE, AP, CGM, SF, NC, YW, MD, JS, MW,
2384 and ME wrote code. AW, and MW developed and deployed the web app. RS, ATA and DP
2385 collected the water samples. DP extracted the samples and acquired the LC-MS/MS data. AKPS,
2386 AW, FO, FR, MN, JB, JJK, EEK, JE, AP, CGM, SF, MRA, TP, NC, MP, CB, BC, AMCR, AC, Fd,
2387 KD, YE, CG, LGG, MH, SH, SK, AK, MCMK, KM, SP, PWP, TS, KSL, PS, ST, GAV, BCW, SX,
2388 MTY, SZ, Md, CB, HPL, CM, JJJvdH, TD, PCD, JS, RS, ATA, ME, and DP tested the protocol,
2389 code and app. CB, JJJvdH, TP, MW, ATA, ME, and DP supervised students and researchers.
2390 MW, AA, ME, and DP supervised the project. AKPS, MN, JB, JJK, EEK, AP, SF, TP, ATA and
2391 DP wrote the manuscript and supplemental information. FO, FR, JE, CGM, MRA, NC, MP, KD,
2392 YE, LGG, MH, SH, PS, GAV, SZ, JJJvdH, TD, TP, PCD, JS, RS, MW, and ME edited and provided
2393 critical feedback on the first draft. All authors edited and approved the final draft.

2394 10. Conflict of Interest

2395 JJJvdH is currently a member of the Scientific Advisory Board of Naicons Srl., Milano, Italy, and
2396 is consulting for Corteva Agriscience, Indianapolis, IN, USA. PCD is a scientific advisor and holds
2397 equity to Cybele and a Co-founder, advisor and holds equity in Ometa, Arome and Enveda with
2398 prior approval by UC-San Diego and consulted in 2023 for DSM animal health. MW is the founder
2399 of Ometa Labs.

2400 11. Additional information & Supplementary information

2401
2402 Supplemental information, including a cheat-sheet, detailed methods for the LC-MS/MS data
2403 acquisition and step-to-step guides for the Python and QIIME2 scripts as well as the web app
2404 are available in the supplemental information.