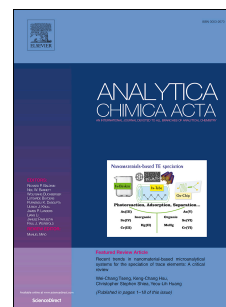# Accepted Manuscript

Chemometric methods in data processing of mass spectrometry-based metabolomics: A review

Lunzhao Yi, Naiping Dong, Yonghuan Yun, Baichuan Deng, Dabing Ren, Shao Liu, Yizeng Liang

Please cite this article as: L. Yi, N. Dong, Y. Yun, B. Deng, D. Ren, S. Liu, Y. Liang, Chemometric methods in data processing of mass spectrometry-based metabolomics: A review, *Analytica Chimica Acta* (2016), doi: 10.1016/j.aca.2016.02.001.

1  Chemometric methods in data processing of mass
2  spectrometry-based metabolomics: A review

3  Lunzhao Yi [a]*, Naiping Dong[c], Yonghuan Yun[b], Baichuan Deng[d], Dabing Ren[a], Shao
4  Liu[e], Yizeng Liang[b]

[a]*Yunnan Food Safety Research Institute, Kunming University of Science and Technology, Kunming, 650500,China*
[b]*College of Chemistry and Chemical Engineering, Central South University, Changsha, 410083, China*
[c]*Department of Applied Biology and Chemical Technology, The Hong Kong Polytechnic University, Hong Kong, 999077, China*
[d]*College of Animal Science, South China Agricultural University, Guangzhou, 510642, China*
[e]*Xiangya hospital, Central South University, Changsha, 410008, China*

5  *Correspondence to: Lunzhao Yi, Yunnan Food safety research institute, Kunming
6  University of Science and Technology, Kunming, 650500, China. Tel.: +86 871
7  65920302. E-mail address: yilunzhao@kmust.edu.cn.

8

## 9  Abstract

10  This review focuses on recent and potential advances in chemometric methods in

11  relation to data processing in metabolomics, especially for data generated from mass

12  spectrometric techniques. Metabolomics is gradually being regarded a valuable and

13  promising biotechnology rather than an ambitious advancement. Herein, we outline

14  significant developments in metabolomics, especially in the combination with modern

15  chemical analysis techniques, and dedicated statistical, and chemometric data

16  analytical strategies. Advanced skills in the preprocessing of raw data, identification

17  of metabolites, variable selection, and modeling are illustrated. We believe that

18  insights from these developments will help narrow the gap between the original

19  dataset and current biological knowledge. We also discuss the limitations and

20  perspectives of extracting information from high-throughput datasets.

21

22  **Keywords:** metabolomics; chemometrics; biomarker; identification of metabolites;
23  data preprocessing; modeling

# Contents

# 1. Introduction

Metabolomics refers to the comprehensive and quantitative analysis of metabolites and aims to gather as much metabolic information as possible from a biological system [1]. It is a reproducible and efficient method that can directly reflect biological events. Metabolomics has recently been upgraded from a promising concept to a widespread and valuable biotechnology. Two modern analytical platforms, namely, nuclear magnetic resonance (NMR) and mass spectrometry (MS), have become the methods of choice for metabolic analysis and are used to generate massive amounts of data to answer various biological questions in metabolomics [2-4].

Improved analytical technologies have gradually caused metabolomics datasets to become larger with more intricate inner structures [5]. Thus, the coverage of

56  metabolomics becomes more comprehensive but will consequently demand more

57  advanced chemometric methods [6]. Metabolomics is either targeted or untargeted. In

58  the targeted approach, specific metabolites of known identity are profiled; good

59  quantitative precision is easily obtained. One disadvantage of this approach, however,

60  is its limitation in terms of the breadth of analysis. The dataset of the target approach

61  is simple. Data analysis often focuses on variable selection and modeling. Untargeted

62  metabolomics aims to simultaneously measure of as many metabolites as possible in a

63  biological specimen. Often, the chemical identities of the MS-resolved peaks are not

64  known a priori, and significant chemical/spectral analysis must be performed to

65  identify the metabolites. Deconvolution and normalization of complex spectra in

66  biological samples is therefore critical for this type of datasets.

67  The raw data from metabolomics presents a gold mine of information [7]. To ensure

68  that the metabolic information is of valuable knowledge, considerable data analysis is

69  required. Chemometrics has become a crucial and dedicated tool for extracting

70  valuable information from data; it presents a complete theory and methodology for

71  every step of metabolomics research, including sampling, experiment design, data

72  pre-processing, metabolite identification, variable selection, and modeling.

73  Chemometrics has thus become one of the cornerstones of metabolomics. However,

74  major changes in the dimensionality and complexity of datasets lead to a significant

75  shift in knowledge discovery. The complexity of metabolomics also presents great

76  challenges on chemometrics to deal with such massive high-dimensional data [6].

77  Several review papers and guide books on metabolomics have been published [8-10],

78  and these works have provided informative and valuable guidance for researchers.

79  Insights into metabolomics experimental skills, including sample preparation and

80  metabolite analysis, have also been revealed [11]. In this review, we describe recent

3

81  advances in chemometric methods for data analysis of metabolomics. This review

82  provides a brief but broad overview of the developed methods, the challenges

83  remaining in the data processing of metabolomics, especially those generated by MS,

84  and perspectives on this topic. Various aspects, including raw data pre-processing,

85  metabolite identification, variable selection and modeling, are discussed. The

86  flowchart of data processing in metabolomics is shown in Figure 1.

87                                    **Insert Figure 1**


88  ## 2. Critique and discussion


89  **2.1 Pre-processing of raw data**

90  Analytical instruments do not provide clean and comparable lists of metabolites. Raw

91  data must be processed to generate a practicable data matrix in a variety of ways [12].

92  The key step is eliminating the variance and bias in the data analysis to reduce the

93  complexity and enhance metabolically significant signals [13]. Consequently, several

94  algorithms have been developed and multiple open source programs have been

95  applied to process raw MS data acquired on liquid chromatography-mass

96  spectrometry (LC-MS) or gas chromatography-mass spectrometry (GC-MS). Among

97  these,   XCMS   (https://xcmsonline.scripps.edu/)   [13,   14],   MZmine

98  (http://sourceforge.net/projects/mzmine/)   [15,   16],   OpenMS

99  (http://open-ms.sourceforge.net/) [17], and MetAlign (http://www.metalign.nl) [18]

100 have attracted particular attentions for their practicability and effectiveness. Most

101 members of the research community of metabolomics work with these tools, and new

102 programs, such as MetSign [19], MSFACTs [20] and MetaboliteDetector [21] have

103 been steadily developed to increase the quality and efficiency of data preprocessing.

104   Most of these tools are freely available. Furthermore, through these tools, the

105   exchange of algorithms and data within the community is convenient. In generally,

106   tools for raw data preprocessing include four basic modules, namely, noise filtering

107   and baseline correction, peak detection and deconvolution, alignment, and

108   normalization. In the following sections, we will introduce different chemometric

109   algorithms and strategies for these modules.

110   2.1.1 Noise filtering and baseline correction

111   Noise filtering is designed to separate component signals from the background

112   originating from the chemical matrix or instrumental interference, remove

113   measurement noise or baseline distortions [9]. Conventionally, during baseline

114   correction of one-way data (e.g. a chromatogram or mass spectrum), the two ends of a

115   signal peak are manually identified by analysts and piecewise linear approximation is

116   then applied to fit a curve as the baseline [22]. However, this procedure is

117   time-consuming, and its accuracy highly depends on the user's operating skills. Thus,

118   numerous algorithms have been developed for better estimation of the baseline. Two

119   powerful algorithms, automatic two-side exponential baseline correction algorithm

120   (ATEB) [23]and adaptive iteratively reweighted penalized least squares (airPLS)

121   [22], were recently developed by Liang's group. These algorithms can automatically

122   and effectively remove the baseline, regardless of whether it is linear or non-linear.

123   Furthermore, unlike methods that require peak detection, these very fast and robust

124   algorithms do not require intervention experience and prior knowledge.

125   For MS-based datasets, the methods for removing random noise are typically

126   implemented by traditional signal processing techniques in chemometrics. Noise

127   filtering of LC-MS data is more complicated than that of GC-MS data because

128   chemical and random noises are both included in the former. Chemical noise is

5

129     typically induced by molecules in buffers and solvents and can be especially strong at

130     the beginning and the end of the elution [24]. This type of noise causes a shift in the

131     baseline in the intermediate mass range of LC-MS spectra. To resolve this problem,

132     several filtering methods have been proposed. For example, Haimi *et al*. fitted the

133     baseline by first segmenting a spectrum and performing linear regression through the

134     lowest points of the smoothed spectrum segments [25]. In addition, baseline removal

135     has also been approached by estimating the background from a two-dimensional

136     intensity image and then removing it with two orthogonal (retention time and *m/z*)

137     one-dimensional passes [26].

138     2.1.2 Peak detection and deconvolution

139     The purpose of peak detection and deconvolution is to identify and quantify the

140     signals corresponding to the molecules (e.g., the metabolites) in a sample [12]. This

141     step is fundamental for downstream data analysis, such as profile alignment or

142     biomarker identification, and can significantly reduce the complexity of the data [9].

143     However, given the complexity of the signals and the multiple sources of noise in data,

144     automatic identification of the noise from compound signals is very difficult. The

145     threshold between noise and a signal is difficult to specify, especially when detecting

146     peaks with low-response values.

147     A peak detection method can identify the true signals correctly and avoid false

148     positives. Unfortunately, high response values do not always guarantee real peaks

149     because some sources of noise can also produce high signals. Conversely, low peaks

150     may correspond to real signals. Therefore, constraints on the peak shapes and criteria

151     of minimal intensity, area or signal-to-noise are widely applied to distinguish real

152     peaks from noise. Several parameters must generally be adjusted to match the

153     characteristics of the MS-based data. Traditionally, peak detection algorithms follow

6

154    two strategies: derivative techniques or matched filter response.

155    Derivative-based peak detection methods make use of the fact that the first derivative

156    of a peak will have a positive-to-negative zero-crossing at the local maxima of a peak

157    [27]. Derivative-based methods commonly require increasingly elaborate

158    pre-processing to prevent compounding noise effects [28, 29]. A slope threshold on is

159    often imposed to avoid false positives.

160    Matched filter methods may become progressively sophisticated as the data

161    complexity increases. One may apply a threshold in the response function to

162    determine the location of chromatographic peaks when applied to chromatographic

163    data by assuming a Gaussian peak shape [30]. A number of popular and open-source

164    software packages, such as XCMS [13] have been developed. XCMS includes three

165    steps: binning, signal determination, and filtering. One weakness of the initially

166    proposed method in XCMS, however, is that the peaks can sometimes be alternatively

167    assigned to two adjacent *m/z* bins. One potential solution to this problem involves

168    combining adjacent extracted ion chromatograms, which represent the analyses of

169    interest. However, this algorithm cannot resolve pairs of co-eluting peaks that fall

170    within half of the *m/z* bin. The developers of XCMS software thus added another

171    algorithm called centWave in later version [31]. The centWave algorithm collects

172    regions containing potentially interesting masses in the raw data and applies

173    continuous wavelet transformation (CWT) and, optionally, Gauss-fitting for

174    chromatographic peak resolution. To circumvent the problems during binning, an

175    alternative fast-computing approach is used in centWave based on the mass accuracy

176    deviation and expected chromatographic peak width. Then, CWT is performed to

177    detect all possible chromatographic peaks. Subsequent filtering is employed to

178    remove candidate peaks in which number of m/z centroids is less then specified

179    threshold. In addition, CWT is also applied to build a robust pattern-matching method

180    for MS peak detection and can be directly used to the raw spectrum. By identifying

181    peaks and assigning a signal-to-noise ratio in the wavelet space according to the

182    two-dimensional CWT coefficient matrix, the pattern matching problem is simplified.

183    Thus, issues surrounding the baseline correction are simultaneously resolved, and the

184    preprocessing steps, such as noise filtering and baseline correction, are not required

185    before peak detection [32].

186    Selecting an optimal threshold for the above mentioned two strategies is a difficult

187    problem but of essential importance that has been thoroughly discussed in various

188    peak detection approaches [27, 33, 34], whereas no general consensus is reached.

189    Some algorithms have recently been developed based on Bayesian inference [35, 36].

190    These algorithms make use of chromatographic information (i.e., the expected width

191    of a single peak and the standard deviation of baseline noise), which is regarded as

192    prior information. Finally, the probability of a signal being a peak is estimated, based

193    on some theories or hypotheses, such as the statistical overlap theory [36].

194    In the high-throughput analysis of metabolites, overlapping peaks are ineluctable.

195    This problem can be resolved by two-dimensional data resolution methods that have

196    been well developed and theorized by the chemometrics community using matrix

197    computation combined with characteristics of spectral data [37-39]. Specifically,

198    multivariate curve resolution-alternating least squares (MCR-ALS) [40] has been

199    extended to processing LC-MS data [41] and shown to be more robust than XCMS

200    [42]. The overlapping peaks can also be resolved by mass spectral deconvolution.

201    Automated mass spectral deconvolution and identification system (AMDIS, NIST)

202    and commercially available tools, such as deconvolution reporting software (DRS,

203    Agilent), AnalyzerPro (SpectralWorks), and ChromaTOF® (LECO), are developed

204  for processing GC-MS data. Most recently, Oliver Fiehn *et al.* [43] proposed an

205  open-source software pipeline, called MS-DIAL, for data-independent acquisition

206  (DIA) - based metabolite identification and quantification by mass spectral

207  deconvolution. MS-DIAL resolves entangled MS/MS spectra by a two-step process:

208  precursor-peak spotting followed by MS/MS-level deconvolution. With this software,

209  DIA can provide high efficacy and accuracy for metabolome coverage.

210  2.1.3 Alignment

211  Alignment of detected features in different samples aims to remove shifts among

212  samples for a given signal to guarantee downstream extraction of useful information.

213  Thus far, several alignment techniques have been developed to minimize run-to-run

214  shifts [44]. To make them applicable to chromatographic systems coupled with

215  sophisticated detection instruments, e.g., LC-MS, which have yielded large amounts

216  of two-dimensional data, the dimensionality must be reduced. The reduction could be

217  achieved by generating integrated peak areas or total ion chromatograms (TICs). For

218  one-dimensional data (such as TICs), some kinds of time alignment procedures could

219  be employed as a useful method for tackling this problem of retention time shifts [45].

220  Examples of these procedures include correlation optimized warping (COW) [46],

221  and dynamic time warping (DTW) [47], recursive alignment by fast Fourier transform

222  (RAFFT) [48]. COW requires large execution times and memory when dealing with

223  huge hyphenated datasets. Artifacts often appear in the fingerprints aligned by DTW

224  because signals are often over-warped when signals are recorded by a mono-channel

225  detector. RAFFT efficiently accelerates the alignment procedure by fast Fourier

226  transform cross-correlation. However, RAFFT may distort the shapes of peaks

227  because it does not consider the peak information when moving segments; this

228  technique only considers the insertion and deletion of data points only at the start and

9

229    end of segments, which may introduce artifacts and remove peak points. Nonlinear

230    retention time shifts often exist for a real sample; thus, a multi-scale peak alignment

231    (MSPA) approach has been proposed. MSPA involves iteratively dividing a

232    chromatogram into smaller segments to solve the problem of nonlinear retention time

233    shifts in alignment. FFT cross correlation is used to estimate candidate shifts and

234    gradually align peaks step by step. A simple example of the application of MSPA

235    method is demonstrated in Figure 2. The retention time shifts of GC-MS TICs in

236    different samples are successfully removed. Other algorithmic alternatives, such as

237    kernel density [13], component-resolving algorithms [49], and progressive clustering

238    [50], among others, exist. Besides, another alignment methods attempt to integrate

239    peak areas. Although time-consuming and meticulous, this approach is considered as

240    the process of "data cleaning" because the retention time shift, noise pollution, and

241    background shift are cleared simultaneously.

242                                    **Insert Figure 2**

243    During dimension reduction, loss of information is inevitable. Addressing this issue

244    involves modeling of the high-dimensional data by multi-way analysis methods,

245    which maintain the so-called two dimensional advantages (e.g., mass spectral

246    information of metabolites). For example, the alignment method by Prakash *et al.* [51]

247    and the ChromAlign method [52] both use the raw high-way data. First, these

248    algorithms construct similarity score matrix for similar spectra between two

249    experimental runs. Dynamic programming is applied to find an optimal path through

250    the matrix and define the mapping of paired spectra. In the method proposed by

251    Pierce *et al.* [53], a piecewise single dimension retention time alignment algorithm is

252    applied to align two-dimensional data. In the continuous profile model (CPM), the

253    two-dimensional data is divided into four *m/z* bins as opposed to the alignment of only

254     a single TIC [54]. In addition, some algorithms align the two-way retention time shift

255     more comprehensively, such as the algorithm using a novel indexing scheme [53].

256     This type of algorithms aligns the fingerprints in different dimensions simultaneously,

257     thereby preserving the separation information in both dimensions.

258     2.1.4 Normalization

259     Normalization removes confounding variations attributed to experimental sources,

260     such as analytical noise or experimental bias, and retains relevant variations attributed

261     to biological events [12]. If the signal of majority of metabolites is stable, simple and

262     efficient normalization could be achieved by calculating the relative ratio of the

263     abundance of analytes to all other peaks, such as the unit norm and median intensities

264     normalization [55]. However, the assumption of negligible overall concentration

265     changes is difficult to satisfy; the total concentrations of analytes may be considerably

266     changed because of laboratory system errors and differences among large scale

267     biological experiments. In this case, scaling based on the total chromatogram may

268     seriously distort the data.

269     Compounds with lower concentrations will be easily altered by analytical noise. To

270     allow the comparison of different metabolites, scaling is required. Autoscaling (1/SD)

271     is the most popular normalization method used in metabolomics; in this method, each

272     variable has equal (unit) variance by multiplying with the inverse of standard

273     deviation (SD). Pareto (1/sqrt(SD)) is softer than autoscaling and can increase the

274     importance of low abundant compounds without significantly amplifying the noise.

275     During data analysis, researchers tend to assume that the total variations originating

276     from sampling, analytical measurements, and biological events are with equal

277     standard deviations and symmetrically around zero [56]. However, this assumption is

278     not satisfied in many cases. Biological effects related to concentration alterations

279 could vary dramatically for different metabolites. Variations related to certain

280 metabolites are considered heteroscedasticity, which could be detrimental to

281 observations of a particular biological situation [56]. A mathematical transformation,

282 such as log transformation [57] or power transformation [58] is helpful to correct the

283 skewed data before modeling. When the relative standard deviation is constant, a log

284 transformation can perfectly remove heteroscedasticity [57]. However, log

285 transformation presents a serious drawback: the transformation approaches minus

286 infinity when the values are transformed as they approach zero. Power transformation

287 does not have the near-zero artifacts and yields results similar to those of log

288 transformation.

289 Another sophisticated strategy for normalization is the internal standards (ISs) method,

290 e.g., isotopically labeled internal standards, and quality control (QC) samples in each

291 data acquisition procedure [59]. Comprehensive and representative IS-based

292 normalization is based on a key assumption that the variance exhibited by ISs solely

293 comes from a component with a systematic error. But, a single IS cannot estimate the

294 systematic error of a complex biological matrix. Multiple ISs work better in this case.

295 Further, IS use must aim to decrease the risk of cross-contribution (CC) which can

296 cause serious loss of information, especially when the interfering analytes are related

297 to the factors of interest in metabolomic datasets. If the masses used for quantifying

298 the IS are carefully selected, this problem can be solved easily [60]. However, this

299 attempt is nontrivial in metabolomics research because the biological sample is too

300 complex. Prediction of which ions will produce cross-interference is difficult.

301 Redestig *et al.* presented an effective normalization algorithm that could compensate

302 for systematic CC effects and improve the normalization of mass spectrometry-based

303 metabolomics data [61]. To image the global variability of a measurement system,

12

304 performing QC before normalization is recommended when visualizing the data by

305 PCA. A QC is a pool of several individuals having similar characteristics. The studied

306 samples are compared with QCs to evaluate their variability. In multivariate statistical

307 analysis, such as PCA, QC samples should appear closely on the scores plot, which

308 indicates that the analytical system has good reproducibility [62].

309 **2.2 Identification of metabolites**

310 Confidently identifying metabolites from MS spectra data has been generally

311 recognized as a significant challenge in the metabolomics community, especially in

312 untargeted analysis, because of the biochemical diversity of metabolites. Given the

313 benefits of advanced computational techniques and methods, advanced mass

314 spectrometry instrumentation, the wealth of knowledge on ion fragmentation, and

315 well-established databases and libraries, especially fruitful works in the past decade,

316 metabolite identification can cover unknowns with reasonable accuracy and could be

317 performed in a high-throughput manner. A variety of overviews have been published

318 on this topic, including basic concepts in compound identification, comprehensive

319 summaries of different identification strategies [63, 64], instructions for practical use

320 [65], and guidelines for beginners of mass spectrometry [66]. Thus, we are going to

321 briefly introduce currently available algorithms and tools valuable for metabolite

322 identification using MS in this section.

323 2.2.1 Metabolite identification using GC-MS

324 GC-MS has been routinely used in metabolomics with mature protocols. Great effort

325 has been made to interpret MS spectra from electron impact (EI) ion sources. The

326 most frequently adopted and reliable method for this is library search, where each

327 experimental MS spectrum is compared with the reference MS spectra in the mass

328　spectral library and the similarity score is calculated for each match. The

329　corresponding library compound gaining the highest similarity score is theoretically

330　considered as the one that generates this experimental spectrum. The commonly

331　adopted mass spectral libraries are listed in Table 1. The main factors that influence

332　search results include the quality of the experimental MS spectra, the size of the mass

333　spectral library, and the similarity score calculation algorithms used [67]. From the

334　arithmetic point of view, the method for calculating the similarity score is the most

335　important factor to consider because the quality of the MS spectra significantly

336　depends on the experiment, and the libraries are generally commercially available and

337　thus cannot be freely configured by users and remain relatively small in size. Previous

338　investigations showed that the most robust similarity score calculation method is the

339　dot product using square-rooted mass spectral intensities [68]. However, no

340　comprehensive comparative investigation is performed for high through-put

341　metabolite identification.

342　**Insert Table 1**

343　Given the complexity of metabolites and their EI-MS spectra, such as the existing of

344　isomers and co-eluted components, a target compound does not ideally gain the

345　highest similarity score but is generally located at a higher rank (e.g., second or third

346　rank, or higher) in the hit list. This approach always requires careful manual checking.

347　Therefore, taking other information, such as the retention index (RI, e.g. Kovat's

348　retention index) of a target compound, into consideration will be very helpful [69, 70].

349　RI is a structurally and physicochemically specific indicator that can effectively

350　differentiate compounds having similar mass spectra. Actually, this indicator and the

351　EI-MS spectrum comprise the widely accepted mass spectral tag (MST) in

352　metabolomics and organize the Golm Metabolome Database(GMD) [71-73] and

14

353 BinBase/FiehnLib [74]. The NIST standard reference database includes a large

354 number of RI values. Another improvement, especially in the case of co-elution, can

355 be achieved by mass spectral deconvolution or two-dimensional data resolution

356 methods (see Section 2.1.2). As GC-MS   instruments with mass analyzers capable of

357 high resolution and accurate mass measurement are now available; the majority of the

358 false matches can also be filtered by considering the accurate masses of the fragments

359 [75].

360 The methods independent of a mass spectral library are to learn the structural features

361 of compounds from their experimental mass spectra and then deduce unknown

362 structures from the features of a given spectrum according to previously constructed

363 learning models. This can be achieved in two ways. The first one involves exhaustion

364 of all possible isomers according to the molecular mass extracted from MS spectra by

365 a structure generation module (e.g., MOLGEN [76] and OMG [77]) and retention of

366 the structures that best explain the spectrum according to fragmentation rules.

367 Machine learning algorithms are generally adopted in this procedure to determine

368 whether a substructure is present in the unknown compound. This step can filter out a

369 large number of isomers that do not contain the identified substructures [78].

370 MOLGEN-MS [79] and MassLib have been developed for this purpose. The

371 web-based algorithm embedded in GMD employs decision trees to predict the 166

372 most common functional groups in metabolites after training known metabolites in

373 GMD with the corresponding mass spectra data and retention indices [80], thereby

374 providing invaluable information for inferring the structures of unknown metabolites.

375 The second approach is based on library search results under the assumption that

376 similar structures have similar spectra. Possible substructures of unknown compounds

377 can be deduced from library compounds with the top similarity scores [81].

15

378     An alternative series of methods directly predict mass spectra for input molecules.

379     Based on the wealth of knowledge on ion fragmentation and aided by advanced

380     computational technologies, accurate prediction of mass spectra has become feasible.

381     Mass Frontier (Thermo Scientific), one of the most commonly adopted software for

382     structure elucidation, uses the HighChem Fragmentation Library, which stores

383     approximately 31,000 fragmentation mechanisms to predict and interpret

384     experimental mass spectra. ACD/MS Fragmenter (ACD/Labs), which is also very

385     powerful for MS spectrum prediction, has gained popularity in the metabolomics

386     community. The freely available tool Mass Spectrum Interpreter, which was released

387     by NIST, uses thermochemical kinetics of general fragmentation reactions

388     summarized from known fragmentation rules to predict mass spectra. Among these

389     powerful methods, a common difficulty is that they cannot effectively extract correct

390     structures from their isomers, as pointed out after comparing different tools [82].

391     However, improvements can be made by combining different tools [83]. In addition to

392     the above methods, by adopting advantages of high resolution GC-MS, unknown

393     compounds can be putatively identified from accurate *m/z* provided by chemical

394     ionization, *in-silico* predicted retention index and fragmentation patterns without

395     requiring any mass spectral library [84, 85]. This trend is analogous to identifying

396     metabolites in high resolution LC-MS, as will be shown below. A practical guide for

397     small molecule structure elucidation with several strategies that differ from above

398     mentioned computational methods can be found in Ref.[86].

399     2.2.2 Metabolite identification using LC-MS

400     For LC-MS, identifying metabolites from MS spectra is not amenable because of the

401     variation of experimental settings, such as chromatographic conditions and mass

402     spectrometry parameters [87]. This step becomes even more serious for discovering

16

403 unknowns from large and complex metabolite space. Additionally, the fragmentation

404 mechanisms during ionization in the LC-MS platform under various activation

405 energies are still unclear. These factors make the confident interpretation of MS

406 spectra derived from different LC-MS and LC-MS$^n$ platforms a significant challenge.

407 Fortunately, recent active studies have made remarkable advances in metabolite

408 identification and several tools and various databases are publically available (see

409 Table 1 and 2). In general, currently available tools are developed based on two

410 aspects of LC-MS data: accurate mass with other information like isotopic

411 distribution and MS/MS spectra.

412 **Insert Table 2**

413 2.2.2.1 Structure inference by accurate mass combined with other information

414 The ability to accurately measure *m/z* is one of the most important features of

415 high-resolution mass spectrometry, which has greatly facilitated the whole MS data

416 analysis workflow. The accurate mass calculated from determined *m/z* is generally the

417 first step [66] because it is the simplest and most straight-forward. The formula

418 generation method or the search of a large compound database or metabolism network

419 can be adopted. For formula generation, all combinations of predefined elements with

420 constraints of element number and mass range are exhausted. A number of tools

421 commercially or freely available have been developed to assist this (see Table 2). As

422 expected, very large number of candidate formulas will be generated, especially for a

423 relatively large molecular mass. This phenomenon makes it impracticable to obtain a

424 single assignment of formula to each *m/z* solely based on the accurate mass. Thus,

425 defining the rules to filter out false positives becomes nontrivial.

426 Among all the developed rules, similarity checking in isotopic distribution is

17

427    commonly accepted as the most critical criterion. Majority of the spurious formulas

428    could be rejected under this checking [88, 89]. Theoretically, each elemental

429    composition or formula has a unique isotopic distribution because different elements

430    have distinct isotopic abundance distributions in nature. Thus, by comparing the

431    instrument-determined isotopic distribution to the simulated one, the formula

432    candidates can be ranked, with the top ones being the most similar via so called

433    spectral comparison [90] or rejected if the relative isotopic abundances (RIA) between

434    the two distributions are unacceptably different. The exploration to precisely simulate

435    isotopic distribution has been undertaken for decades and several tools are now freely

436    available [91]. If the resolution of an MS instrument is high enough, formulas can be

437    exclusively identified from the RIA of a single element. This strategy is now extended

438    and confirmed with higher-resolution instruments for high-throughput metabolomics

439    analysis [92]. However, high RIA measurement errors can appear in peaks with a low

440    signal-to-noise ratio (S/N), low *m/z*, and the presence of co-eluting species [93-96].

441    These factors will terribly mislead the identification results [94]. Unfortunately, the

442    systematic evaluation of the influence of RIA measurement error on formulae

443    inference is not performed. A suggestion for eliminating this influence can be setting

444    a larger error tolerance during comparison [95]. Whereas cautions still should be

445    proceeded with when using RIA to identify metabolites and additional information is

446    required.

447    The second rule is to check whether the generated formulas are reasonable as

448    candidates of metabolites. The famous "Seven Golden Rules" was defined after

449    statistically analyzing formulas extracted from Wiley and NIST02 mass spectral

450    database and the Dictionary of Natural Products [88] and has been demonstrated to be

451    an efficient tool in metabolomics. An updated version of these rules is defined

18

452    recently after analyzing large scale formulas in the PubChem database [97].

453    Once formulas are determined or ranked, decoding them to known metabolites in

454    LC-MS feature annotation is subsequently performed, typically by searching large

455    chemical substance databases [98, 99]. The databases frequently adopted in

456    metabolomics are listed in Table 1. Further annotation of ion species can be realized

457    by prior biological knowledge from lists of expected metabolites of the analyzed

458    organism. Metabolites in biological samples are biochemically connected (e.g.,

459    chemical transformation) rather than randomly mixed [100]. Thus, the metabolite

460    candidates are mapped onto metabolism networks to gain confident identification

461    [101-103]. For example, MI-Pack maps mass spectral peaks onto the KEGG network

462    database [104] and uses the rigidly defined mass error surface of mass differences

463    between substrate-product pairs derived from the database for metabolite

464    identification [103]. Significant reduction of false negatives and false positives is

465    consequently obtained. This approach is advantageous for metabolite identification

466    and mining related subnetworks, which represent the activity or functions of the

467    metabolites, as demonstrated in recent works [105, 106].

468    Besides mapping ions to molecular databases, mining relationships between extracted

469    ion features to annotate these ions has also been proven to be a highly effective

470    strategy. This approach can be executed because LC-MS can detect ion series (so

471    called satellite ions) of a metabolite generated by fragmentation reactions during

472    ionization, including neutral losses and ions with different adducts [107, 108]. This

473    process can generate an *in silico* ion network that reveals relationships between

474    metabolites, also known as metabolic biotransformation [100]. CAMERA [109],

475    IDEOM [110], and MAIT [111] *etc*. were developed in this manner.

19

476　2.2.2.2 Metabolite identification by $MS^n$

477　$MS^n$ is a highly effective technique for structure elucidation. As an indispensable part

478　of the LC-MS system, ionized molecules or molecules in the *m/z* range specified by

479　instruments are gradually dissociated into charged or neutral pieces by hard ionization

480　methods such as the collision-induced dissociation (CID). Recording all the charged

481　fragments and precursor ion forms the $MS^n$ spectrum. This $MS^n$ spectrum generation

482　procedure demonstrates that a molecule's structure can be readily deduced from its

483　$MS^n$ spectrum. Moreover, strategies for interpreting GC-MS spectra (e.g., library

484　search or mass spectrum prediction) can be applied in this deduction. Therefore,

485　several $MS^n$ spectral libraries and computational methods for spectral prediction or

486　structure elucidation are developed (Table1 and 2). The experimental conditions (e.g.,

487　collision energy) in $MS^n$ analysis are not as standardized as in GC-MS analysis.

488　Furthermore, the sizes of currently constructed libraries are much smaller compared

489　with the whole metabolism or structure databases and other factors [112, 113]. Thus,

490　metabolite identification via spectral library search is not as popular in $MS^n$ analysis

491　as in GC-MS analysis. Consequently, much more studies are focused on developing

492　computational methods to interpret $MS^n$ spectra without querying spectral libraries.

493　The algorithms employed in currently developed software for computational $MS^n$ can

494　be categorized into three basic approaches, namely, mass spectrum prediction, *in*

495　*silico* fragmentation, and *de novo* elucidation [114]. Mass spectrum prediction, which

496　is mainly applied for $MS^2$, has been well studied in EI spectrum interpretation. This

497　process is also a basic and highly important module in peptide identification under

498　hypothesis-driven proteomics. The enormous diversity of small compounds continues

499　to considerably challenge accurate $MS^2$ spectral prediction. To predict the $MS^2$

500　spectrum for a given structure, Mass Frontier extracts all possible reactions that can

20

501 occur during the fragmentation of this structure from its own fragmentation reaction

502 library to generate rules for the prediction of fragments and intensities. ACD/MS

503 Fragmenter handles spectrum prediction in a similar way. MetISIS uses a

504 machine-learning algorithm to learn CID kinetics from lipid experimental $MS^2$ spectra

505 to predict lipid spectra *in silico* [115]. A fragment ion prediction algorithm embedded

506 in MyCompoundID website (http://www.mycompoundid.org/) adopts a "chopping"

507 program to predict the bond cleavage of metabolites to generate theoretical $MS^2$

508 spectra for database search [116]. Instead of directly predicting mass spectra, *in silico*

509 fragmentation attempts to elucidate a structure from all candidates that best explains

510 the given $MS^2$ spectrum. This approach was first employed in EPIC using a bond

511 disconnection algorithm to exhaust all possible substructures of a molecule and

512 compare the substructures to formulas inferred from fragment ions. Then relevant

513 structures were listed for user confirmation [117]. Later, FiD [118] and

514 Mass-MetaSite [119] were developed on the basis of bond dissociation mechanism,

515 and MetFrag extended this procedure [120] by considering rearrangement reactions

516 during molecule fragmentation. An alternative procedure was implemented in

517 FingerID by calculating the likelihood between metabolites in a database and a given

518 experimental $MS^2$ spectrum in a feature space called fingerprints using an support

519 vector machine (SVM) model [121]. This model was obtained by training fingerprints

520 extracted from the Mass Bank MS/MS ($MS^2$) spectral library. CFM calculated the

521 likelihood between database metabolites and given $MS^2$ spectra in accordance with

522 the competitive fragmentation process learned from a spectral library using the

523 expectation maximum algorithm [122].

524 *De novo* analysis, however, infers structures from the observed fragments in a given

525 $MS^n$ spectrum. This approach first determines the formulas of fragments according to

526   their high resolution *m/z* and then deduces the structure of a precursor ion using these

527   formulas and the known fragmentation pathways that generate these ions. To date, the

528   most appropriate method employed for this deduction appears to be the construction

529   of a fragmentation tree with nodes being fragment formulas, edges being neutral

530   losses, and the root being the precursor [123, 124]. Therefore, with an appropriate

531   scoring scheme, an experimental $MS^n$ spectrum can be identified by extracting the

532   most optimal fragmentation tree defined by the scores. Even so, the later portion of

533   this procedure has been demonstrated to be extremely computationally intensive,

534   despite already attaining the precursor formulas [125]. This obstacle can be partly

535   solved by heuristic methods [126] and several tools, such as $SIRIUS^2$ [123, 124] and

536   MAGMa [127].

### 2.3 Variable selection

538   Variable selection aims to extracting important metabolites from a mass of

539   metabolites detected by mass spectrometry that can help us to answer biological

540   questions at hand, which plays an essential role in metabolomics. From statistical

541   point of view, this is an optimization approach that discovers an optimal variable

542   combination from the considerable body of variables. However, this process faces a

543   great challenge to address the NP-hard problem called "large p, small n problem"

544   [128]. To date, numerous variable selection methods specific to this problem have

545   been proposed. Some of these suggested strategies are based on statistical features of

546   variables, whereas some are based on the optimization algorithm. Herein, we divide

547   these methods into two kinds of approaches as follows: variable ranking and variable

548   subset selection [129].

549

550 2.3.1Variable ranking

551 Variable ranking is mostly used in revealing informative metabolites or biomarkers.

552 The process of ranking assigns a measure of importance to each variable on the basis

553 of certain criteria. Many PLS-based criteria are frequently employed for variable

554 ranking[130], including PLS loading weights (LW) [131], variable importance on

555 projection (VIP) scores [132], regression coefficient (RC) [133], target projection (TP)

556 [134], and selectivity ratio (SR) [135]. To date, VIP is the most popular one in

557 metabolomics. Yi *et al.* [136] reported that VIP exhibited better efficiency than LW

558 and RC for the metabolomics dataset of nasopharyngeal carcinoma patients. However,

559 for another dataset, the comparison result between different variable ranking methods

560 might be different [137, 138]. Because the efficiency of these methods is

561 data-dependent, it is hard to say which one is the best. We should know that various

562 variable ranking methods are most likely to generate different variable ranking results

563 due to their different principles. Recently, Yun *et al.* use rank aggregation method to

564 emerge all different ranking lists into a final aggregated variable ranking list for

565 biomarker discovery [139]. It is a good attempt to handle this problem. In addition,

566 variable ranking can be conducted based on statistical features between variables and

567 classification label.

568 2.3.2 Variable subset selection

569 Subset selection refers to the search for an optimal subset from all variables that

570 satisfy an optimality criterion. Any variable ranking method can be transformed into a

571 variable subset selection algorithm by introducing a threshold on the variable

572 importance values. The assignment of this threshold can be subjective or achieved by

573 statistical method [129]. Usually, a trade-off between model prediction accuracy and

574 the number of selected variables is considered. The most straightforward proposal for

23

575　this purpose is to use cross validation (CV) procedure to determine the threshold. This

576　approach estimates the generalization error using different number of variables and

577　chooses the number that minimizes the prediction error (CV error). That is, after

578　ranking variables from the most important to the least by some criteria (e.g., VIP),

579　models are built by adding these variables sequentially until all are included, and CV

580　error obtained by each model is recorded. The best variable subset can then be

581　determined to be the first $n$ variables if minimum CV error is achieved after adding

582　$n$th variable. In addition, some criteria related to the classification algorithm can be

583　employed for subset selection. The objective function is a pattern classifier, which

584　evaluates variable subsets according to their predictive accuracy by statistical

585　re-sampling (e.g., bootstrapping) or CV. Usually, optimization algorithm is combined

586　with the classification algorithm. And, variable subset selection seeks the optimal or

587　near-optimal subset with respect to an objective function. For example, genetic

588　algorithm - Bayesian network （ GA-BayesN ） approach [140] combines the

589　optimization algorithm GA with a classifier. Compared with the variable ranking

590　method, subset selection generally achieves better prediction accuracy because the

591　latter considers the specific interactions between the classifier and dataset. In the

592　process, subset selection utilizes a mechanism to avoid overfitting through

593　re-sampling or CV measures of prediction accuracy. However, the approach entails

594　training of a classifier for each variable subset, leading to low execution and high

595　computation. Moreover, the solution lacks generality because subset selection

596　combines the bias of the classifier with the fitness evaluation function.

597　2.3.3 Variable selection considering the interaction effect among variables

598　In fact, finding an optimal subset or variable ranking is not always preferred unless

599　the interaction among multiple variables is considered. The collective effect of

600　variables should be considered because the joint performance of a set of variables is

601　better than the additive independent contributions of its individuals [141]. To address

602　this problem, Zhao and Liu introduced a variable subset selection method, called

603　INTERACT [142]. This approach is based on inconsistency and symmetrical

604　uncertainty measurements for finding interacting features. The group proposed

605　variable interactions can be implicitly managed with a carefully designed variable

606　evaluation metric and a search strategy with a specially designed data structure. The

607　metric and strategy together take the combination effects of variables into

608　considerations when performing variable selection. The method proposed in

609　Breiman's work [143] somewhat considers the combination effects of variables on the

610　basis of the random forest (RF) and permutation test. The variable importance is

611　assessed by the percent increase of misclassification error when the variable is

612　randomly permuted in a RF. However, all variables are involved in the RF model,

613　thus, providing a good reflection of the synergetic effect among multiple variables is

614　difficult to accomplish.

615　Recently, Liang's group proposed a new strategy for variable selection, called model

616　population analysis (MPA) [144]. This method provides a general framework for the

617　development of data analysis methods. Figure 3 illustrates the outline of the MPA.

618　MPA involves three steps. Firstly, (1) sampling method (e.g., Monte Carlo sampling

619　(MCS)) is employed to randomly produce N sub-datasets (e.g., 10,000). Then, (2) a

620　sub-model is built on each sub-dataset. And finally, (3) statistical analysis is

621　employed to evaluate outcomes of interest (e.g., prediction errors) for all established

622　N sub-models. With this approach, the variables are identified as informative,

623　uninformative, or interfering variables according to the differences between the cases

25

624    and control samples. Figure 4 illustrates the prediction error distributions of the three

625    kinds of variables after permutation. Uninformative and interfering variables are

626    useless because of their potential undesirable influence on the modeling. Thus,

627    discovering the optimal variable subset or ranking in the informative variables can

628    produce compelling results.

629                                    **Insert Figure 3**

630                                    **Insert Figure 4**

631    Subwindow permutation analysis (SPA) [145] combines the above-mentioned

632    concepts on the MCS method and MPA. SPA assesses each variable's importance on

633    the basis of the sub-models obtained by MCS technique. Informative variables are

634    identified and ranked by $p$ values obtained by the Mann-Whitney U test on two

635    distributions of prediction errors. Another method, margin influence analysis (MIA)

636    [146] is also based on the concepts of MCS and MPA. Although designed to operate

637    with SVM in identifying informative variables, MIA also offers a measure for each

638    variable on the basis of the differences between the prediction errors from the

639    inclusion and exclusion of this variable. However, the chance of each variable to be

640    sampled by MCS is not the same. Some variables are selected more frequently than

641    others; hence, assessing the importance of each variable using the above introduced

642    strategy does not appear appropriate. So, a new sampling method in the variable space,

643    called binary matrix sampling (BMS) [147], was proposed. This method not only

644    considers the synergetic effect among multiple variables, but also guarantees that each

645    variable is selected with equal probability and a population of different variable

646    combinations is concurrently generated. With this population of variable subset, Yun

647    *et al.* introduced a method called variable importance analysis, which is based on

648    random variable combination (VIAVC) [137]. VIAVC employs the MPA strategy and

26

649 finds the optimal subset of variables by observing the differences between the

650 prediction errors of inclusion and exclusion of each variable. Meanwhile, Deng *et al*.

651 developed an optimization algorithm called variable iterative space shrinkage

652 approach (VISSA) to determine optimal variable combinations [148]. Each variable is

653 assigned a weight according to its importance during modeling in VISSA. The weight

654 of each variable accumulates through an iterative procedure and the variables are

655 selected when their weights reach "1". Two rules are highlighted in the VISSA

656 algorithm. First, the variable space shrinks smoothly in each step. Second, the variable

657 space is optimized in each step.

658 Although the above mentioned methods considered the synergetic effect among

659 multiple variables, these approaches rarely investigate the complementary information

660 between variables. By contrast, the variable complementary network (VCN) is an

661 overall method that visualizes the complementary processes among biological

662 variables [149]. VCN accumulates the information from several classification models

663 obtained by MCS in variable space, quantitatively computes the complementary

664 information between variables. Thus, it can effectively discover biomarkers with the

665 aid of mutual associations among metabolites.   For comparison, Table 3 lists several

666 variable selection methods.

667 **Insert Table 3**

668 **2.4 Modeling of the data**

669 To explore the high-dimensional metabolomics datasets and discover valuable

670 information on biological events, a number of machine-learning methods have been

671 applied. Main characteristics of the machine learning methods which will be

672 described below are summarized in Table 4. It contains the category, advantages and

27

673   disadvantages of each method, and also some applications in metabolomics.

674                                  **Insert Table 4**

675   2.4.1 Unsupervised methods

676   Unsupervised methods are usually used to explore the overall structure of a dataset,

677   finding trends and groupings within the dataset. These methods contribute an

678   unbiased view of the data. Several unsupervised methods are available, during which

679   principal component analysis (PCA), hierarchical cluster analysis (HCA), and

680   self-organization mapping (SOM) are the most frequently used examples in

681   metabolomics.

682   PCA transforms the high-dimensional variables into a small number of orthogonal

683   factors, called PCs, containing the largest variance [150, 151]. PCA provides the

684   projection of samples into low dimensional (usually two- or three-dimensional) PC

685   space, enabling the visualization of the sample distribution. HCA aims to group

686   relatively similar samples in one cluster and relatively dissimilar objects in another

687   [152, 153].

688   SOM is a neural-network algorithm [154]. For high-dimensional data, SOM can form

689   a non-linear projection on a regular, low-dimensional grid. The clustering in the data

690   space and the metric-topological relations of the data items is clearly visible. SOM is

691   a useful tool to characterize metabolic patterns and interrelationship between samples

692   [155-157]. For example, similar responses on primary and secondary metabolites

693   were characterized in microorganisms across stimuli using MS-based metabolomics

28

694    and SOM [155].

695    PARAFAC2 [158] is an extension of PARAFAC [159] which can be used to model

696    three-way data with a trilinear structure. PARAFAC2 can be considered as the

697    generalization of PCA to a higher order of data. It allows the simultaneous processing

698    of all samples, deconvolution of metabolites, elimination of chromatographic baseline,

699    and alignment of retention time shifts. PARAFAC2 can provide simple and robust

700    models upon the application of some constraints. The advantage of PARAFAC2 is its

701    ability of finding and modeling the shifted peaks of the same chemical compounds,

702    with the disadvantage of being sensitive to noise [160, 161]. Goodacre *et al.* [162]

703    employed PARAFAC2 to model the metabolic profiles of meat and characterise the

704    hygiene status of pork chops which undergo a spoilage process.

705    2.4.2 Supervised methods

706    Supervised techniques support a priori known data structures to train patterns and

707    rules to predict new data, which can be classified as linear methods, such as PLS-DA,

708    linear discriminant analysis (LDA), orthogonal projections to latent structures

709    discriminant analysis (OPLS-DA), and non-linear methods, including RF and SVM

710    *etc.*

711    LDA attempts to find a linear function on the basis of original variables, which

712    maximizes the ratio of between-class variance and minimizes the ratio of within-class

713    variance [152]. LDA is a fast and powerful tool for discriminant analysis, in which

714    parameter optimization is not necessary. The number of samples must be larger than

715    that of the variables, ensuring that the inverse of the covariance matrix can be

29

716    obtained [163].

717    The most widely used supervised method for classification in chemometrics is

718    PLS-DA [164], which is a combination of PLS regression and LDA. One advantage

719    of PLS-DA is its ability to handle highly collinear data. Moreover, PLS-DA can

720    provide excellent insights into the cause of discrimination by checking the behavior of

721    variables (e.g. variable importance, see Section 2.3.1). As such, PLS-DA is also a

722    useful tool in biomarker discovery. The recent modification of PLS-DA is the

723    OPLS-DA [165]. The systematic variations in data matrix X can be split into two

724    parts through the orthogonal signal correction (OSC) technique [166]: one part

725    exhibits linear responsiveness, whereas another is linearly orthogonal to the response.

726    OPLS-DA supposes that only the variance related to the response is useful for

727    modeling [167]. It gives better visualization and interpretation than PLS-DA [168]

728    and has been widely applied in modeling and biomarker discovery in metabolomics

729    [169-171].

730     2.4.3 Non-linear methods

731    Complex interactions occur in different levels of biological organizations; hence,

732    biological processes commonly follow a non-linear response. In these cases,

733    non-linear pattern recognition methods are required to characterize metabolomics data.

734    Many non-linear techniques have been proposed in pattern recognition and machine

735    learning research fields. Among these methods, kernel-PLS, RF and SVMs are three

736    popular methods used in metabolomics.

737    Kernel-based models transform data using some specific functions called kernels. By

30

738  using the kernel transformation, researchers can transform the non-linear problem of

739  the original data into a higher-dimensional feature space. Afterward, the non-linear

740  problem becomes linear and can be solved easily. The kernel functions appear in

741  various types, and users can choose appropriate kernel transformation for a certain

742  dataset. Positive semi-definite is one requirement of the kernel matrix [172].

743  Meanwhile, the dot product is the simplest kernel function for the data matrix. The

744  radial basic function is another frequently used kernel function that requires tuning of

745  parameters relating to the width of the Gaussian. Kernel-based classification methods,

746  such as kernel Fisher discriminant analysis (K-FDA) [173], kernel PLS (K-PLS) [174],

747  and kernel OPLS (KO-PLS) [175] have been developed and all exhibit obvious

748  advantages in solving non-linear problems.

749  SVM is another powerful kernel-based classifier that utilizes a set of objects called

750  support vectors to define decision boundaries and separate binary class[176]. SVM

751  focuses on finding a hyper-plane that splits two classes perfectly, whereas the

752  thickness of the margins is maximized. Hence, for each class, the distance of the plane

753  to the data point is the closest [177, 178]. If a point is situated on the wrong side of

754  the margin, the margin is maximized by penalizing the point. The step can split the

755  overlapping classes. Support vectors are the points on the boundary or on the wrong

756  side of the margin supporting the split. When classes are separated by a non-linear

757  boundary, the kernel method is used to find the boundary. SVM is particularly suitable

758  for the data of small sample sizes. The scheme is also capable of handling both linear

759  and non-linear problems of classification by applying linear and non-linear kernels.

31

760 The major disadvantage of SVM is that the model is lack of transparency and variable

761 importance is difficult to obtain. Another disadvantage is that it does not provide a

762 universal means of solving non-linear problems. Hence, kernel functions should be

763 selected discreetly [179]. It has been applied in toxicology research [180, 181], food

764 research [182], and *etc*.

765 RF [143] is an ensemble-learning method that consists of a large number of

766 classification and regression trees (CART). It is highly powerful classifier for

767 high-dimensional data. A random resampling method with replacement called

768 bootstrapping [183] is used to select training samples from the original samples

769 (bootstrap samples) to construct a classification tree. Bootstrapping is carried out

770 many times to build a large group of simple CARTs. Model accuracy is improved with

771 the help of bootstrapping using resample means to estimate sample means [184]. Two

772 powerful and efficient machine-learning techniques, bagging and random feature

773 selection are employed in RF. For bagging, each CART is trained on the bootstrap

774 samples of the training dataset. Predictions are obtained from the majority of votes of

775 the CARTs. During RF model construction, only about two-thirds of training samples

776 are used due to the intrinsic property of bootstrap sampling. Thus the remaining

777 samples can serve as an internal testing set to monitor the prediction error termed

778 out-of-bag error (OOB error). Besides, variable importance can also be obtained by

779 comparing OOB error difference between normal variable and its random permutation,

780 as has been introduced in previous section. RF has shown better performance than

781 many of the classifiers such as PLS-DA and OPLS-DA with external validation [185].

32

782    Other examples also showed that RF and other approaches could be the alternatives to

783    PLS-DA [186]. It has been applied to metabolomics research of hepatocellular

784    carcinoma [187], breast cancer [188], metabolic syndrome [189], and *etc*.

785    2.4.4 Model tuning and model validation

786    The tuning of parameters is of great importance when building a model. CV [190] is

787    the most commonly used model tuning method because it selects a model on the basis

788    of prediction ability. Leave-one-out CV, K-fold CV [191] and Monte Carlo CV [192]

789    are important branches of CV. Recently, CV has faced up some criticisms. For

790    example, it may provide exceedingly optimistic results of the model prediction ability

791    [193]. An alternative is to use double CV (DCV) which involves two loops: the inner

792    loop is used for model tuning, and outer loop is adopted for model validation [194].

793    Model validation is a process on deciding whether results quantify hypothesized

794    relationships between variables and responses and provide accurate estimation of the

795    model prediction ability. Supervised machine-learning methods, such as PLS-DA,

796    hold a high tendency for over-fitting, especially in high dimensional data [195, 196].

797    Thus, a careful model validation is desired.

798    Several criteria can evaluate the prediction ability of a model including sensitivity,

799    specificity, accuracy, the receiver operating characteristic (ROC) curve, and the

800    cross-validated coefficient of determination ($Q^2$). For a perfect classification, the

801    value of specificity should be close to 1, and 1- specificity should be preferably close

802    to 0. When the area under the ROC curve (AUC) is closer to 1, the method performs

803    better. Recently, a criterion was developed by combing $Q^2$ and model stability (S)

33

804　[197]. The results show that, when a clear maximum of $Q^2$ is not obtained, S can

805　provide additional information of over-fitting and it helps in finding the optimal nLVs.

806　We believe that the criterion will be efficient for model selection of metabolomics.

807　The most common strategies and recommended for model validation are independent

808　test set, CV and permutation test. Ideally, model validation employs an independent

809　test set assumed to be representative and independent from the training data. A

810　number of algorithms can be adopted to divide samples into training and test sets,

811　including the Duplex algorithm [198], Kennard-and-Stone algorithm [199], and SPXY

812　algorithm [200]. However, the ideal situation is usually unsatisfied in actual settings,

813　often resulting in bias findings.

814　In CV, model tuning and model validation processes are carried out simultaneously.

815　When the optimal model parameter is determined, the characteristics of prediction

816　ability, such as $Q^2$, are obtained by tuning parameters. However, in DCV the model

817　tuning and model validation processes are carried out separately using inner loop and

818　outer loop, respectively. DCV has shown more accurate estimations of error rates than

819　six-fold CV [194].

820　Permutation test is another powerful approach for model validation. The class labels

821　of samples are permutated randomly in a permutation test. By repeating the

822　permutation test numerous times, a group of "wrong" models are built, and the

823　distribution for accuracy, $Q^2$, and AUC can be obtained. For a validated model, the

824　difference between the "right" models and the "wrong" models should be significant.

825　This difference can be characterized by statistical hypothesis testing. The permutation

34

826    test also offers many applications in metabolomics studies [201, 202].

827    The modeling of metabolomics data is a kind of systematic work. For exploratory

828    studies, unsupervised methods, such as PCA, provide an informative first look at the

829    dataset structures and relationships between groups. Then, supervised methods, such

830    as PLS-DA and OPLS-DA, are applied to classify the samples as well as discover

831    biomarkers. When these classifiers fail to work properly, non-linear models SVM and

832    RF are applied to further explore the non-linear relationship within the data. In

833    addition, the parameters of each model should be well tuned and the model should be

834    validated with caution to ensure its prediction ability for future samples.

835    **2.5. One eye on the future**

836    To date, numerous authors have demonstrated that data processing based on an

837    individual datasets limit the complete understanding of the chemical complexity of

838    the metabolome. Substantial data and information is generated from numerous

839    experimental platforms (e.g., NMR, GC-MS or LC-MS). Consequently, the

840    combination of information becomes increasingly necessary and important in

841    extending metabolite coverage and characterizing biological systems [5]. The greatest

842    future challenge is on how to efficiently integrate massive information from various

843    sources (i.e., data fusion problem). Merging information from multiple datasets with

844    different structural characteristics and extracting the common or distinctive features

845    will unquestionably form a crucial element for the more comprehensive prospect of

846    metabolomics.

847  An increasing number of papers have been published to discuss the problem of data

848  fusion since 2005 [5, 56, 203]. Data fusions often focus on handling multiple datasets

849  generated by several analytical platforms and analyzing longitudinal metabolomic

850  data with time-resolved models. Boccard and Rudaz proposed the four main

851  approaches to data fusion: low-level, mid-level, high-level and kernel-based data

852  fusion [5]. Low-level fusion simply merges data matrices from different platforms

853  into a single matrix for regression analysis or discriminant analysis. Mid-level fusion

854  first extracts relevant features from each data source and then concatenates these

855  features into a single matrix. In high-level fusion, separate models are obtained from

856  each data source and the results of each model are combined to obtain the final

857  decision [204]. Kernel-based methods employ kernel functions to transform data into

858  high-dimensional feature spaces and generate kernel matrices. The kernel matrices are

859  then merged to construct a single matrix for modeling [205]. The selection of data

860  fusion methods depends on the difference between data sources. Data sources with

861  larger differences often entail higher levels of data fusion. So far, methods for data

862  fusion mainly focus on the low and middle levels [203, 206, 207]. Further fusion

863  includes the integration with various "omics" fields, such as genomics,

864  transcriptomics, and proteomics. These methods are all effective strategies for

865  describing a whole biological system. However, we should be careful to avoid the

866  network discordance when metabolomics are integrated with other "omics" [208].

867 ## 3. Conclusions

868 In summary, metabolomics plays an essential role in basic research for elucidating

869 environmental effects, gene functions, and defining cellular processes. To date,

870 research on this field entails much exercise of caution with regard to data acquisition,

871 processing, and information interpretation because of the numerous limitations related

872 to data processing in metabolomics. We herein emphasize four issues, which are of

873 great importance for data processing in metabolomics as follows. 1) Automatic and

874 effective data preprocessing remains a difficult task, especially for the detection,

875 alignment, and deconvolution of peaks with low responses. 2) The confident

876 identification of unknown metabolites from complex MS spectra data remains as a

877 great challenge. 3) NP-hard problems in variable selection must be addressed but

878 barely solved by all researchers. 4) New efficient model validation methods and

879 indices are urgently desired. Furthermore, these methods must be carefully selected in

880 practice to guarantee that the objective models are fully validated and with good

881 prediction ability for future actual samples. All of these problems, along with the

882 high-dimensional characteristics of metabolomic datasets pose numerous fundamental

883 questions in chemometrics. Chemometrics is facing enormous challenges to develop

884 robust and efficient methods to answer various biological questions derived from

885 metabolomics. We believe that this review can guide practitioners of metabolomics,

886 and provide insights into its present uses as well as new data processing applications.

887 ## Conflicts of interest statement

888 The author declares no conflicts of interest.

37

## Acknowledgements

## References

[1] R. Goodacre, S. Vaidyanathan, W.B. Dunn, G.G. Harrigan, D.B. Kell, Metabolomics by numbers: acquiring and understanding global metabolite data, Trends in biotechnology, 22 (2004) 245-252.

[2] H.K. Kim, Y.H. Choi, R. Verpoorte, NMR-based plant metabolomics: where do we stand, where do we go?, Trends in biotechnology, 29 (2011) 267-275.

[3] J.W. Allwood, R. Goodacre, An introduction to liquid chromatography–mass spectrometry instrumentation applied in plant metabolomic analyses, Phytochemical analysis, 21 (2010) 33-47.

[4] L. Yi, C. Song, Z. Hu, L. Yang, L. Xiao, B. Yi, W. Jiang, Y. Cao, L. Sun, A metabolic discrimination model for nasopharyngeal carcinoma and its potential role in the therapeutic evaluation of radiotherapy, Metabolomics, 10 (2014) 697-708.

[5] J. Boccard, S. Rudaz, Harnessing the complexity of metabolomic data with chemometrics, Journal of Chemometrics, 28 (2014) 1-9.

[6] J. van der Greef, A.K. Smilde, Symbiosis of chemometrics and metabolomics: past, present, and future, Journal of Chemometrics, 19 (2005) 376-386.

[7] R. Goodacre, Making sense of the metabolome using evolutionary computation: seeing the wood with the trees, Journal of experimental botany, 56 (2005) 245-254.

[8] A.H. BaniMustafa, N.W. Hardy, A Strategy for Selecting Data Mining Techniques in Metabolomics, Plant Metabolomics, Springer2012, pp. 317-333.

[9] M. Katajamaa, M. Orešič, Data processing for mass spectrometry-based metabolomics, Journal of Chromatography A, 1158 (2007) 318-328.

[10] A.M. De Livera, M. Sysi-Aho, L. Jacob, J.A. Gagnon-Bartsch, S. Castillo, J.A. Simpson, T.P. Speed, Statistical Methods for Handling Unwanted Variation in Metabolomics Data, Analytical chemistry, 87 (2015) 3606-3615.

[11] M. Ernst, D.B. Silva, R.R. Silva, R.Z. Vêncio, N.P. Lopes, Mass spectrometry in plant metabolomics strategies: from analytical platforms to data acquisition and processing, Natural product reports, (2014).

[12] S. Castillo, P. Gopalacharyulu, L. Yetukuri, M. Orešič, Algorithms and tools for the preprocessing of LC–MS metabolomics data, Chemometrics and Intelligent Laboratory Systems, 108 (2011) 23-32.

[13] C.A. Smith, E.J. Want, G. O'Maille, R. Abagyan, G. Siuzdak, XCMS: processing mass spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and identification, Analytical chemistry, 78 (2006) 779-787.

[14] H. Benton, D. Wong, S. Trauger, G. Siuzdak, XCMS2: processing tandem mass spectrometry data

for metabolite identification and structural characterization, Analytical chemistry, 80 (2008) 6382-6389.

[15] M. Katajamaa, J. Miettinen, M. Orešič, MZmine: toolbox for processing and visualization of mass spectrometry based molecular profile data, Bioinformatics, 22 (2006) 634-636.

[16] T. Pluskal, S. Castillo, A. Villar-Briones, M. Orešič, MZmine 2: modular framework for processing, visualizing, and analyzing mass spectrometry-based molecular profile data, BMC bioinformatics, 11 (2010) 395.

[17] M. Sturm, A. Bertsch, C. Gropl, A. Hildebrandt, R. Hussong, E. Lange, N. Pfeifer, O. Schulz-Trieglaff, A. Zerck, K. Reinert, O. Kohlbacher, OpenMS-An open-source software framework for mass spectrometry, BMC Bioinformatics, 9 (2008).

[18] R.C. De Vos, S. Moco, A. Lommen, J.J. Keurentjes, R.J. Bino, R.D. Hall, Untargeted large-scale plant metabolomics using liquid chromatography coupled to mass spectrometry, Nature protocols, 2 (2007) 778-791.

[19] X. Wei, W. Sun, X. Shi, I. Koo, B. Wang, J. Zhang, X. Yin, Y. Tang, B. Bogdanov, S. Kim, MetSign: A computational platform for high-resolution mass spectrometry-based metabolomics, Analytical chemistry, 83 (2011) 7668-7675.

[20] A.L. Duran, J. Yang, L. Wang, L.W. Sumner, Metabolomics spectral formatting, alignment and conversion tools (MSFACTs), Bioinformatics, 19 (2003) 2283-2293.

[21] K. Hiller, J. Hangebrauk, C. Jäger, J. Spura, K. Schreiber, D. Schomburg, MetaboliteDetector: comprehensive analysis tool for targeted and nontargeted GC/MS based metabolome analysis, Analytical chemistry, 81 (2009) 3429-3439.

[22] Z.-M. Zhang, S. Chen, Y.-Z. Liang, Baseline correction using adaptive iteratively reweighted penalized least squares, Analyst, 135 (2010) 1138-1146.

[23] X. Liu, Z. Zhang, Y. Liang, P.F. Sousa, Y. Yun, L. Yu, Baseline correction of high resolution spectral profile data based on exponential smoothing, Chemometrics and Intelligent Laboratory Systems, 139 (2014) 97-108.

[24] M. Hilario, A. Kalousis, C. Pellegrini, M. Mueller, Processing and classification of protein mass spectra, Mass spectrometry reviews, 25 (2006) 409-449.

[25] P. Haimi, A. Uphoff, M. Hermansson, P. Somerharju, Software tools for analysis of mass spectrometric lipidome data, Analytical chemistry, 78 (2006) 8324-8331.

[26] M. Bellew, M. Coram, M. Fitzgibbon, M. Igra, T. Randolph, P. Wang, D. May, J. Eng, R. Fang, C. Lin, A suite of algorithms for the comprehensive analysis of complex protein mixtures using high-resolution LC-MS, Bioinformatics, 22 (2006) 1902-1909.

[27] G. Vivó-Truyols, J. Torres-Lapasió, A. Van Nederkassel, Y. Vander Heyden, D. Massart, Automatic program for peak detection and deconvolution of multi-overlapped chromatographic signals: Part I: Peak detection, Journal of Chromatography A, 1096 (2005) 133-145.

[28] K.M. Pierce, R.E. Mohler, A Review of chemometrics applied to comprehensive two-dimensional separations from 2008–2010, Separation & Purification Reviews, 41 (2012) 143-168.

[29] S. Krishnan, J.T. Vogels, L. Coulier, R.C. Bas, M.W. Hendriks, T. Hankemeier, U. Thissen, Instrument and process independent binning and baseline correction methods for liquid chromatography–high resolution-mass spectrometry deconvolution, Analytica Chimica Acta, 740 (2012) 12-19.

[30] R. Danielsson, D. Bylund, K.E. Markides, Matched filtering with background suppression for improved quality of base peak chromatograms and mass spectra in liquid chromatography–mass spectrometry, Analytica Chimica Acta, 454 (2002) 167-184.

972  [31] R. Tautenhahn, C. Böttcher, S. Neumann, Highly sensitive feature detection for high resolution
973  LC/MS, BMC bioinformatics, 9 (2008) 504.

974  [32] P. Du, W.A. Kibbe, S.M. Lin, Improved peak detection in mass spectrum by incorporating
975  continuous wavelet transform-based pattern matching, Bioinformatics, 22 (2006) 2059-2065.

976  [33] K.C. Leptos, D.A. Sarracino, J.D. Jaffe, B. Krastins, G.M. Church, MapQuant: Open‐source
977  software for large‐scale protein quantification, Proteomics, 6 (2006) 1770-1782.

978  [34] C.A. Hastings, S.M. Norton, S. Roy, New algorithms for processing and peak detection in liquid
979  chromatography/mass spectrometry data, Rapid communications in mass spectrometry, 16 (2002)
980  462-467.

981  [35] G. Vivó-Truyols, Bayesian approach for peak detection in two-dimensional chromatography,
982  Analytical chemistry, 84 (2012) 2622-2630.

983  [36] M. Lopatka, G. Vivó-Truyols, M. Sjerps, Probabilistic peak detection for first-order
984  chromatographic data, Analytica Chimica Acta, 817 (2014) 9-16.

985  [37] Y.Z. Liang, O.M. Kvalheim, Resolution of two-way data: theoretical background and practical
986  problem-solving - Part 1: Theoretical background and methodology, Fresen J Anal Chem, 370 (2001)
987  694-704.

988  [38] L.W. Hantao, H.G. Aleme, M.P. Pedroso, G.P. Sabin, R.J. Poppi, F. Augusto, Multivariate curve
989  resolution combined with gas chromatography to enhance analytical separation in complex samples: A
990  review, Analytica Chimica Acta, 731 (2012) 11-23.

991  [39] C. Ruckebusch, L. Blanchet, Multivariate curve resolution: A review of advanced and tailored
992  applications and challenges, Analytica Chimica Acta, 765 (2013) 28-36.

993  [40] R. Tauler, Multivariate curve resolution applied to second order data, Chemometrics and
994  Intelligent Laboratory Systems, 30 (1995) 133-146.

995  [41] E. Gorrochategui, J. Jaumot, R. Tauler, A protocol for LC-MS metabolomic data processing using
996  chemometric tools, Protocol Exchange, (2015).

997  [42] M. Navarro-Reig, J. Jaumot, A. Garcia-Reiriz, R. Tauler, Evaluation of changes induced in rice
998  metabolome by Cd and Cu exposure using LC-MS with XCMS and MCR-ALS data analysis strategies,
999  Analytical and Bioanalytical Chemistry, 407 (2015) 8835-8847.

1000  [43] H. Tsugawa, T. Cajka, T. Kind, Y. Ma, B. Higgins, K. Ikeda, M. Kanazawa, J. VanderGheynst, O. Fiehn,
1001  M. Arita, MS-DIAL: data-independent MS/MS deconvolution for comprehensive metabolome analysis,
1002  Nature methods, 12 (2015) 523-526.

1003  [44] R. Smith, D. Ventura, J.T. Prince, LC-MS alignment in theory and practice: a comprehensive
1004  algorithmic review, Briefings in bioinformatics, 16 (2015) 104-117.

1005  [45] K.J. Johnson, B.W. Wright, K.H. Jarman, R.E. Synovec, High-speed peak matching algorithm for
1006  retention time alignment of gas chromatographic data for chemometric analysis, Journal of
1007  Chromatography A, 996 (2003) 141-155.

1008  [46] N.-P.V. Nielsen, J.M. Carstensen, J. Smedsgaard, Aligning of single and multiple wavelength
1009  chromatographic profiles for chemometric data analysis using correlation optimised warping, Journal
1010  of Chromatography A, 805 (1998) 17-35.

1011  [47] V. Pravdova, B. Walczak, D. Massart, A comparison of two algorithms for warping of analytical
1012  signals, Analytica Chimica Acta, 456 (2002) 77-92.

1013  [48] J.W. Wong, C. Durante, H.M. Cartwright, Application of fast Fourier transform cross-correlation for
1014  the alignment of large chromatographic and spectral datasets, Analytical chemistry, 77 (2005)
1015  5655-5661.

[49] V.P. Andreev, T. Rejtar, H.-S. Chen, E.V. Moskovets, A.R. Ivanov, B.L. Karger, A universal denoising and peak picking algorithm for LC-MS based on matched filtration in the chromatographic time domain, Analytical chemistry, 75 (2003) 6314-6326.

[50] D.P. De Souza, E.C. Saunders, M.J. McConville, V.A. Likić, Progressive peak clustering in GC-MS Metabolomic experiments applied to Leishmania parasites, Bioinformatics, 22 (2006) 1391-1396.

[51] A. Prakash, P. Mallick, J. Whiteaker, H. Zhang, A. Paulovich, M. Flory, H. Lee, R. Aebersold, B. Schwikowski, Signal maps for mass spectrometry-based comparative proteomics, Molecular & cellular proteomics, 5 (2006) 423-432.

[52] R.G. Sadygov, F. Martin Maroto, A.F. Hühmer, ChromAlign: a two-step algorithmic procedure for time alignment of three-dimensional LC-MS chromatographic surfaces, Analytical chemistry, 78 (2006) 8207-8217.

[53] K.M. Pierce, L.F. Wood, B.W. Wright, R.E. Synovec, A comprehensive two-dimensional retention time alignment algorithm to enhance chemometric analysis of comprehensive two-dimensional separation data, Analytical chemistry, 77 (2005) 7735-7743.

[54] J. Listgarten, R.M. Neal, S.T. Roweis, P. Wong, A. Emili, Difference detection in LC-MS data for protein biomarker discovery, Bioinformatics, 23 (2007) e198-e204.

[55] W. Wang, H. Zhou, H. Lin, S. Roy, T.A. Shaler, L.R. Hill, S. Norton, P. Kumar, M. Anderle, C.H. Becker, Quantification of proteins and metabolites by mass spectrometry without isotopic labeling or spiked standards, Analytical chemistry, 75 (2003) 4818-4826.

[56] R.A. van den Berg, H.C. Hoefsloot, J.A. Westerhuis, A.K. Smilde, M.J. van der Werf, Centering, scaling, and transformations: improving the biological information content of metabolomics data, BMC genomics, 7 (2006) 142.

[57] O.M. Kvalheim, F. Brakstad, Y. Liang, Preprocessing of analytical profiles in the presence of homoscedastic or heteroscedastic noise, Analytical chemistry, 66 (1994) 43-51.

[58] R. Sokal, F. Rohlf, Assumptions of analysis of variance, Biometry: The Principles and Practice of Statistics in Biological Research. 3rd ed. New York: WH Freeman, (1995) 396-406.

[59] H.G. Gika, G. Theodoridis, J. Extance, A.M. Edge, I.D. Wilson, High temperature-ultra performance liquid chromatography–mass spectrometry for the metabonomic analysis of Zucker rat urine, Journal of Chromatography B, 871 (2008) 279-287.

[60] R. Liu, D. Lin, W. Chang, C. Liu, W. Tsay, J. Li, T. Kuo, Issues to address when isotopically labeled analogues of analytes are used as internal standards, Anal. Chem, 74 (2002) 618AJ626A.

[61] H. Redestig, A. Fukushima, H. Stenlund, T. Moritz, M. Arita, K. Saito, M. Kusano, Compensation for Systematic Cross-Contribution Improves Normalization of Mass Spectrometry Based Metabolomics Data, Analytical chemistry, 81 (2009) 7974-7980.

[62] H.G. Gika, E. Macpherson, G.A. Theodoridis, I.D. Wilson, Evaluation of the repeatability of ultra-performance liquid chromatography–TOF-MS for global metabolic profiling of human urine samples, Journal of Chromatography B, 871 (2008) 299-305.

[63] D.S. Wishart, Computational strategies for metabolite identification in metabolomics, Bioanalysis, 1 (2009) 1579-1596.

[64] T. Kind, O. Fiehn, Advances in structure elucidation of small molecules using mass spectrometry, Bioanal Rev, 2 (2010) 23-60.

[65] D.G. Watson, A rough guide to metabolite identification using high resolution liquid chromatography mass spectrometry in metabolomic profiling in metazoans, Comput Struct Biotechnol J, 4 (2013) e201301005.

1060 [66] M. Holcapek, R. Jirasko, M. Lisa, Basic rules for the interpretation of atmospheric pressure
1061 ionization mass spectra of small molecules, J Chromatogr A, 1217 (2010) 3908-3921.

1062 [67] I. Koo, S. Kim, X. Zhang, Comparative analysis of mass spectral matching-based compound
1063 identification in gas chromatography-mass spectrometry, J Chromatogr A, 1298 (2013) 132-138.

1064 [68] S.E. Stein, D.R. Scott, Optimization and testing of mass spectral library search algorithms for
1065 compound identification, J Am Soc Mass Spectrom, 5 (1994) 859-866.

1066 [69] W.B. Dunn, D. Broadhurst, P. Begley, E. Zelena, S. Francis-McIntyre, N. Anderson, M. Brown, J.D.
1067 Knowles, A. Halsall, J.N. Haselden, A.W. Nicholls, I.D. Wilson, D.B. Kell, R. Goodacre, Procedures for
1068 large-scale metabolic profiling of serum and plasma using gas chromatography and liquid
1069 chromatography coupled to mass spectrometry, Nat Protoc, 6 (2011) 1060-1083.

1070 [70] J. Kopka, Current challenges and developments in GC-MS based metabolite profiling technology,
1071 Journal of Biotechnology, 124 (2006) 312-322.

1072 [71] J. Kopka, N. Schauer, S. Krueger, C. Birkemeyer, B. Usadel, E. Bergmuller, P. Dormann, W.
1073 Weckwerth, Y. Gibon, M. Stitt, L. Willmitzer, A.R. Fernie, D. Steinhauser, GMD@CSB.DB: the Golm
1074 Metabolome Database, Bioinformatics, 21 (2005) 1635-1638.

1075 [72] C. Wagner, M. Sefkow, J. Kopka, Construction and application of a mass spectral and retention
1076 time index database generated from plant GC/EI-TOF-MS metabolite profiles, Phytochemistry, 62
1077 (2003) 887-900.

1078 [73] N. Schauer, D. Steinhauser, S. Strelkov, D. Schomburg, G. Allison, T. Moritz, K. Lundgren, U.
1079 Roessner-Tunali, M.G. Forbes, L. Willmitzer, A.R. Fernie, J. Kopka, GC-MS libraries for the rapid
1080 identification of metabolites in complex biological samples, Febs Letters, 579 (2005) 1332-1337.

1081 [74] T. Kind, G. Wohlgemuth, D.Y. Lee, Y. Lu, M. Palazoglu, S. Shahbaz, O. Fiehn, FiehnLib: Mass Spectral
1082 and Retention Index Libraries for Metabolomics Based on Quadrupole and Time-of-Flight Gas
1083 Chromatography/Mass Spectrometry, Anal Chem, 81 (2009) 10038-10048.

1084 [75] N.W. Kwiecien, D.J. Bailey, M.J.P. Rushp, J.S. Cole, A. Ulbrich, A.S. Hebert, M.S. Westphall, J.J. Coon,
1085 High-Resolution Filtering for Improved Small Molecule Identification via GC/MS, Anal Chem, 87 (2015)
1086 8328-8335.

1087 [76] C. Benecke, R. Grund, R. Hohberger, A. Kerber, R. Laue, T. Wieland, Molgen(+), a Generator of
1088 Connectivity Isomers and Stereoisomers for Molecular-Structure Elucidation, Analytica Chimica Acta,
1089 314 (1995) 141-147.

1090 [77] J.E. Peironcely, M. Rojas-Cherto, D. Fichera, T. Reijmers, L. Coulier, J.L. Faulon, T. Hankemeier, OMG:
1091 Open Molecule Generator, J Cheminform, 4 (2012) 21.

1092 [78] E.L. Schymanski, C. Meinert, M. Meringer, W. Brack, The use of MS classifiers and structure
1093 generation to assist in the identification of unknowns in effect-directed analysis, Analytica Chimica
1094 Acta, 615 (2008) 136-147.

1095 [79] A. Kerber, R. Laue, M. Meringer, K. Varmuza, MOLGEN-MS: Evaluation of low resolution electron
1096 impact mass spectra with MS classification and exhaustive structure generation, in: E. Gelpi (Ed.)
1097 Advances in Mass Spectrometry 15, Wiley, 2001, pp. 939-940.

1098 [80] J. Hummel, N. Strehmel, J. Selbig, D. Walther, J. Kopka, Decision tree supported substructure
1099 prediction of metabolites from GC-MS profiles, Metabolomics, 6 (2010) 322-333.

1100 [81] S.E. Stein, Chemical substructure identification by mass spectral library searching, J Am Soc Mass
1101 Spectrom, 6 (1995) 644-655.

1102 [82] E.L. Schymanski, M. Meringer, W. Brack, Matching Structures to Mass Spectra Using
1103 Fragmentation Patterns: Are the Results As Good As They Look?, Anal Chem, 81 (2009) 3608-3617.

42

[83] E.L. Schymanski, C.M.J. Gallampois, M. Krauss, M. Meringer, S. Neumann, T. Schulze, S. Wolf, W. Brack, Consensus Structure Elucidation Combining GC/EI-MS, Structure Generation, and Calculated Properties, Anal Chem, 84 (2012) 3287-3295.

[84] S. Kumari, D. Stevens, T. Kind, C. Denkert, O. Fiehn, Applying In-Silico Retention Index and Mass Spectra Matching for Identification of Unknown Metabolites in Accurate Mass GC-TOF Mass Spectrometry, Anal Chem, 83 (2011) 5895-5902.

[85] O. Fiehn, J. Kopka, R.N. Trethewey, L. Willmitzer, Identification of uncommon plant metabolites based on calculation of elemental compositions using gas chromatography and quadrupole mass spectrometry, Anal Chem, 72 (2000) 3573-3580.

[86] L.X. Zhang, C.L. Tang, D.S. Cao, Y.X. Zeng, B.B. Tan, M.M. Zeng, W. Fan, H.B. Xiao, Y.Z. Liang, Strategies for structure elucidation of small molecules using gas chromatography-mass spectrometric data, Trac-Trend Anal Chem, 47 (2013) 37-46.

[87] J.M. Halket, D. Waterman, A.M. Przyborowska, R.K.P. Patel, P.D. Fraser, P.M. Bramley, Chemical derivatization and mass spectral libraries in metabolic profiling by GC/MS and LC/MS/MS, Journal of Experimental Botany, 56 (2005) 219-243.

[88] T. Kind, O. Fiehn, Seven Golden Rules for heuristic filtering of molecular formulas obtained by accurate mass spectrometry, BMC Bioinformatics, 8 (2007) 105.

[89] J.C.L. Erve, M. Gu, Y.D. Wang, W. DeMaio, R.E. Talaat, Spectral Accuracy of Molecular Ions in an LTQ/Orbitrap Mass Spectrometer and Implications for Elemental Composition Determination, J Am Soc Mass Spectr, 20 (2009) 2058-2069.

[90] Y.D. Wang, M. Cu, The Concept of Spectral Accuracy for MS, Anal Chem, 82 (2010) 7055-7062.

[91] D. Valkenborg, I. Mertens, F. Lemiere, E. Witters, T. Burzykowski, The isotopic distribution conundrum, Mass Spectrometry Reviews, 31 (2012) 96-109.

[92] T. Nagao, D. Yukihira, Y. Fujimura, K. Saito, K. Takahashi, D. Miura, H. Wariishi, Power of isotopic fine structure for unambiguous determination of metabolite elemental compositions: in silico evaluation and metabolomic application, Anal Chim Acta, 813 (2014) 70-76.

[93] Y. Xu, J.F. Heilier, G. Madalinski, E. Genin, E. Ezan, J.C. Tabet, C. Junot, Evaluation of Accurate Mass and Relative Isotopic Abundance Measurements in the LTQ-Orbitrap Mass Spectrometer for Further Metabolomics Database Building, Anal Chem, 82 (2010) 5490-5501.

[94] B.P. Koch, T. Dittmar, M. Witt, G. Kattner, Fundamentals of molecular formula assignment to ultrahigh resolution mass data of natural organic matter, Anal Chem, 79 (2007) 1758-1763.

[95] R.J.M. Weber, A.D. Southam, U. Sommer, M.R. Viant, Characterization of Isotopic Abundance Measurements in High Resolution FT-ICR and Orbitrap Mass Spectra for Improved Confidence of Metabolite Identification, Anal Chem, 83 (2011) 3737-3743.

[96] A. Knolhoff, J. Callahan, T. Croley, Mass Accuracy and Isotopic Abundance Measurements for HR-MS Instrumentation: Capabilities for Non-Targeted Analyses, J Am Soc Mass Spectr, 25 (2014) 1285-1294.

[97] A. Lommen, Ultrafast PubChem Searching Combined with Improved Filtering Rules for Elemental Composition Analysis, Anal Chem, 86 (2014) 5463-5469.

[98] Z.J. Zhu, A.W. Schultz, J.H. Wang, C.H. Johnson, S.M. Yannone, G.J. Patti, G. Siuzdak, Liquid chromatography quadrupole time-of-flight mass spectrometry characterization of metabolites guided by the METLIN database, Nature Protocols, 8 (2013) 451-460.

[99] J. Little, A. Williams, A. Pshenichnov, V. Tkachenko, Identification of "Known Unknowns" Utilizing Accurate Mass Data and ChemSpider, J Am Soc Mass Spectr, 23 (2012) 179-185.

43

1148 [100] R. Breitling, S. Ritchie, D. Goodenowe, M.L. Stewart, M.P. Barrett, Ab initio prediction of
1149 metabolic networks using Fourier transform mass spectrometry data, Metabolomics, 2 (2006)
1150 155-164.

1151 [101] G.T. Gipson, K.S. Tatsuoka, B.A. Sokhansanj, R.J. Ball, S.C. Connor, Assignment of MS-based
1152 metabolomic datasets via compound interaction pair mapping, Metabolomics, 4 (2008) 94-103.

1153 [102] S. Rogers, R.A. Scheltema, M. Girolami, R. Breitling, Probabilistic assignment of formulas to mass
1154 peaks in metabolomics experiments, Bioinformatics, 25 (2009) 512-518.

1155 [103] R.J.M. Weber, M.R. Viant, MI-Pack: Increased confidence of metabolite identification in mass
1156 spectra by integrating accurate masses and metabolic pathways, Chemometr Intell Lab, 104 (2010)
1157 75-82.

1158 [104] H. Ogata, S. Goto, K. Sato, W. Fujibuchi, H. Bono, M. Kanehisa, KEGG: Kyoto Encyclopedia of
1159 Genes and Genomes, Nucleic Acids Res, 27 (1999) 29-34.

1160 [105] H. Doerfler, X. Sun, L. Wang, D. Engelmeier, D. Lyon, W. Weckwerth,
1161 mzGroupAnalyzer--predicting pathways and novel chemical structures from untargeted
1162 high-throughput metabolomics data, PLoS One, 9 (2014) e96188.

1163 [106] S. Li, Y. Park, S. Duraisingham, F.H. Strobel, N. Khan, Q.A. Soltow, D.P. Jones, B. Pulendran,
1164 Predicting Network Activity from High Throughput Metabolomics, PLoS Comput Biol, 9 (2013)
1165 e1003123.

1166 [107] N. Huang, M.M. Siegel, G.H. Kruppa, F.H. Laukien, Automation of a Fourier transform ion
1167 cyclotron resonance mass spectrometer for acquisition, analysis, and E-mailing of high-resolution
1168 exact-mass electrospray ionization mass spectral data, J Am Soc Mass Spectr, 10 (1999) 1166-1173.

1169 [108] M. Brown, W.B. Dunn, P. Dobson, Y. Patel, C.L. Winder, S. Francis-McIntyre, P. Begley, K. Carroll, D.
1170 Broadhurst, A. Tseng, N. Swainston, I. Spasic, R. Goodacre, D.B. Kell, Mass spectrometry tools and
1171 metabolite-specific databases for molecular identification in metabolomics, Analyst, 134 (2009)
1172 1322-1332.

1173 [109] C. Kuhl, R. Tautenhahn, C. Bottcher, T.R. Larson, S. Neumann, CAMERA: an integrated strategy for
1174 compound spectra extraction and annotation of liquid chromatography/mass spectrometry data sets,
1175 Anal Chem, 84 (2012) 283-289.

1176 [110] D.J. Creek, A. Jankevics, K.E. Burgess, R. Breitling, M.P. Barrett, IDEOM: an Excel interface for
1177 analysis of LC-MS-based metabolomics data, Bioinformatics, 28 (2012) 1048-1049.

1178 [111] F. Fernandez-Albert, R. Llorach, C. Andres-Lacueva, A. Perera, An R package to analyse LC/MS
1179 metabolomic data: MAIT (Metabolite Automatic Identification Toolkit), Bioinformatics, 30 (2014)
1180 1937-1939.

1181 [112] S. Stein, Mass Spectral Reference Libraries: An Ever-Expanding Resource for Chemical
1182 Identification, Anal Chem, 84 (2012) 7274-7282.

1183 [113] E. Werner, J.F. Heilier, C. Ducruix, E. Ezan, C. Junot, J.C. Tabet, Mass spectrometry for the
1184 identification of the discriminating signals from metabolomics: Current status and future trends, J
1185 Chromatogr B, 871 (2008) 143-163.

1186 [114] F. Hufsky, K. Scheubert, S. Bocker, Computational mass spectrometry for small-molecule
1187 fragmentation, Trac-Trend Anal Chem, 53 (2014) 41-48.

1188 [115] L.J. Kangas, T.O. Metz, G. Isaac, B.T. Schrom, B. Ginovska-Pangovska, L. Wang, L. Tan, R.R. Lewis,
1189 J.H. Miller, In silico identification software (ISIS): a machine learning approach to tandem mass
1190 spectral identification of lipids, Bioinformatics, 28 (2012) 1705-1713.

1191 [116] T. Huan, C.Q. Tang, R.H. Li, Y. Shi, G.H. Lin, L. Li, MyCompoundID MS/MS Search: Metabolite

1192    Identification Using a Library of Predicted Fragment-Ion-Spectra of 383,830 Possible Human

1193    Metabolites, Analytical Chemistry, 87 (2015) 10619-10626.

1194    [117] A.W. Hill, R.J. Mortishire-Smith, Automated assignment of high-resolution collisionally activated

1195    dissociation mass spectra using a systematic bond disconnection approach, Rapid Commun Mass Sp,

1196    19 (2005) 3111-3118.

1197    [118] M. Heinonen, A. Rantanen, T. Mielikainen, J. Kokkonen, J. Kiuru, R.A. Ketola, J. Rousu, FiD: a

1198    software for ab initio structural identification of product ions from tandem mass spectrometric data,

1199    Rapid Commun Mass Spectrom, 22 (2008) 3043-3052.

1200    [119] B. Bonn, C. Leandersson, F. Fontaine, I. Zamora, Enhanced metabolite identification with MS(E)

1201    and a semi-automated software for structural elucidation, Rapid Commun Mass Spectrom, 24 (2010)

1202    3127-3138.

1203    [120] S. Wolf, S. Schmidt, M. Muller-Hannemann, S. Neumann, In silico fragmentation for computer

1204    assisted identification of metabolite mass spectra, BMC Bioinformatics, 11 (2010) 148.

1205    [121] M. Heinonen, H. Shen, N. Zamboni, J. Rousu, Metabolite identification and molecular fingerprint

1206    prediction through machine learning, Bioinformatics, 28 (2012) 2333-2341.

1207    [122] F. Allen, R. Greiner, D. Wishart, Competitive fragmentation modeling of ESI-MS/MS spectra for

1208    putative metabolite identification, Metabolomics, (2014) 1-13.

1209    [123] S. Bocker, F. Rasche, Towards de novo identification of metabolites by analyzing tandem mass

1210    spectra, Bioinformatics, 24 (2008) i49-i55.

1211    [124] F. Rasche, A. Svatos, R.K. Maddula, C. Bottcher, S. Bocker, Computing fragmentation trees from

1212    tandem mass spectrometry data, Anal Chem, 83 (2011) 1243-1251.

1213    [125] F. Hufsky, M. Rempt, F. Rasche, G. Pohnert, S. Böcker, De novo analysis of electron impact mass

1214    spectra using fragmentation trees, Analytica Chimica Acta, 739 (2012) 67-76.

1215    [126] I. Rauf, F. Rasche, F. Nicolas, S. Böcker, Finding Maximum Colorful Subtrees in Practice, in: B. Chor

1216    (Ed.) Lect N Bioinformat, Springer Berlin Heidelberg2012, pp. 213-223.

1217    [127] L. Ridder, J.J.J. van der Hooft, S. Verhoeven, R.C.H. de Vos, R. van Schaik, J. Vervoort,

1218    Substructure-based annotation of high-resolution multistage MSn spectral trees, Rapid

1219    Communications in Mass Spectrometry, 26 (2012) 2461-2471.

1220    [128] J. Boccard, J.L. Veuthey, S. Rudaz, Knowledge discovery in metabolomics: an overview of MS data

1221    handling, Journal of separation science, 33 (2010) 290-304.

1222    [129] I. Narsky, F.C. Porter, Methods for Variable Ranking and Selection,  Statistical Analysis

1223    Techniques in Particle Physics, Wiley-VCH Verlag GmbH & Co. KGaA2013, pp. 385-415.

1224    [130] T. Mehmood, K.H. Liland, L. Snipen, S. Sæbø, A review of variable selection methods in partial

1225    least squares regression, Chemometrics and Intelligent Laboratory Systems, 118 (2012) 62-69.

1226    [131] S. Wold, M. Sjöström, L. Eriksson, Partial Least Squares Projections to Latent Structures (PLS) in

1227    Chemistry,   Encyclopedia of Computational Chemistry, John Wiley & Sons, Ltd2002.

1228    [132] S. Favilla, C. Durante, M.L. Vigni, M. Cocchi, Assessing feature relevance in NPLS models by VIP,

1229    Chemom. Intell. Lab. Syst, 129 (2013) 76-86.

1230    [133] S. Wold, M. Sjöström, L. Eriksson, PLS-regression: a basic tool of chemometrics, Chemom. Intell.

1231    Lab. Syst, 58 (2001) 109-130.

1232    [134] T. Rajalahti, R. Arneberg, A.C. Kroksveen, M. Berle, K.-M. Myhr, O.M. Kvalheim, Discriminating

1233    Variable Test and Selectivity Ratio Plot: Quantitative Tools for Interpretation and Variable (Biomarker)

1234    Selection in Complex Spectral or Chromatographic Profiles, Analytical Chemistry, 81 (2009) 2581-2590.

1235    [135] O.M. Kvalheim, Interpretation of partial least squares regression models by means of target

45

1236     projection and selectivity ratio plots, J. Chemometr, 24 (2010) 496-504.

1237     [136] L.Z. Yi, N.P. Dong, S.T. Shi, B.C. Deng, Y.H. Yun, Z.B. Yi, Y. Zhang, Metabolomic identification of

1238     novel biomarkers of nasopharyngeal carcinoma, Rsc Advances, 4 (2014) 59094-59101.

1239     [137] Y.-H. Yun, F. Liang, B.-C. Deng, G.-B. Lai, C.M.V. Goncalves, H.-M. Lu, J. Yan, X. Huang, L.-Z. Yi, Y.-Z.

1240     Liang, Informative metabolites identification by variable importance analysis based on random

1241     variable combination, Metabolomics, 11 (2015) 1539-1551.

1242     [138] M. Farrés, S. Platikanov, S. Tsakovski, R. Tauler, Comparison of the variable importance in

1243     projection (VIP) and of the selectivity ratio (SR) methods for variable selection and interpretation,

1244     Journal of Chemometrics, 29 (2015) 528-536.

1245     [139] Y.-H. Yun, B.-C. Deng, D.-S. Cao, W.-T. Wang, Y.-Z. Liang, Variable importance analysis based on

1246     rank aggregation with applications in metabolomics for biomarker discovery, Analytica Chimica Acta,

1247     (2016).

1248     [140] E. Correa, R. Goodacre, A genetic algorithm-Bayesian network approach for the analysis of

1249     metabolomics and spectroscopic data: application to the rapid identification of Bacillus spores and

1250     classification of Bacillus species, BMC bioinformatics, 12 (2011) 33.

1251     [141] D. Anastassiou, Computational analysis of the synergy among multiple interacting genes,

1252     Molecular Systems Biology, 3 (2007) n/a-n/a.

1253     [142] Z. Zhao, H. Liu, Searching for interacting features in subset selection, Intelligent Data Analysis, 13

1254     (2009) 207-228.

1255     [143] L. Breiman, Random forests, Machine Learning, 45 (2001) 5-32.

1256     [144] H.-D. Li, Y.-Z. Liang, Q.-S. Xu, D.-S. Cao, Model population analysis for variable selection, J.

1257     Chemometr, 24 (2010) 418-423.

1258     [145] H.-D. Li, M.-M. Zeng, B.-B. Tan, Y.-Z. Liang, Q.-S. Xu, D.-S. Cao, Recipe for revealing informative

1259     metabolites based on model population analysis, Metabolomics, 6 (2010) 353-361.

1260     [146] H.-D. Li, Y.-Z. Liang, Q.-S. Xu, D.-S. Cao, Recipe for Uncovering Predictive Genes Using Support

1261     Vector Machines Based on Model Population Analysis, IEEE. ACM. T. Comput. Bi, 8 (2011) 1633-1641.

1262     [147] H. Zhang, H. Wang, Z. Dai, M.-s. Chen, Z. Yuan, Improving accuracy for cancer classification with

1263     a new algorithm for genes selection, BMC Bioinformatics, 13 (2012) 1-20.

1264     [148] B.-c. Deng, Y.-h. Yun, Y.-z. Liang, L.-z. Yi, A novel variable selection approach that iteratively

1265     optimizes variable space using weighted binary matrix sampling, Analyst, 139 (2014) 4836-4845.

1266     [149] H.-D. Li, Q.-S. Xu, W. Zhang, Y.-Z. Liang, Variable complementary network: a novel approach for

1267     identifying biomarkers and their mutual associations, Metabolomics, 8 (2012) 1218-1226.

1268     [150] J.E. Jackson, A User's Guide to Principal Components, Wiley, New York, 1991.

1269     [151] J. Xu, F.L. Hu, W. Wang, X.C. Wan, G.H. Bao, Investigation on biochemical compositional changes

1270     during the microbial fermentation process of Fu brick tea by LC-MS based metabolomics, Food Chem,

1271     186 (2015) 176-184.

1272     [152] A.R. Webb, Statistical pattern recognition, John Wiley & Sons2003.

1273     [153] L. Jing, Z.T. Lei, G.W. Zhang, A.C. Pilon, D.V. Huhman, R.J. Xie, W.P. Xi, Z.Q. Zhou, L.W. Sumner,

1274     Metabolite profiles of essential oils in citrus peels and their taxonomic implications, Metabolomics, 11

1275     (2015) 952-963.

1276     [154] T. Kohonen, S.-O. Maps, Springer series in information sciences, Self-organizing maps, 30 (1995).

1277     [155] C.R. Goodwin, B.C. Covington, D.K. Derewacz, C.R. McNees, J.P. Wikswo, J.A. McLean, B.O.

1278     Bachmann, Structuring Microbial Metabolic Responses to Multiplexed Stimuli via Self-Organizing

1279     Metabolomics Maps, Chemistry & biology, (2015).

46

[156] J.K. Kim, M.R. Cho, H.J. Baek, T.H. Ryu, C.Y. Yu, M.J. Kim, E. Fukusaki, A. Kobayashi, Analysis of metabolite profile data using batch-learning self-organizing maps, Journal of Plant Biology, 50 (2007) 517-521.

[157] A.D. Patterson, H. Li, G.S. Eichler, K.W. Krausz, J.N. Weinstein, A.J. Fornace, F.J. Gonzalez, J.R. Idle, UPLC-ESI-TOFMS-Based Metabolomics and Gene Expression Dynamics Inspector Self-Organizing Metabolomic Maps as Tools for Understanding the Cellular Response to Ionizing Radiation, Analytical Chemistry, 80 (2008) 665-674.

[158] J.M. Amigo, T. Skov, R. Bro, J. Coello, S. Maspoch, Solving GC-MS problems with parafac2, TrAC Trends in Analytical Chemistry, 27 (2008) 714-725.

[159] R. Bro, PARAFAC. Tutorial and applications, Chemometr Intell Lab, 38 (1997) 149-171.

[160] B. Khakimov, J.M. Amigo, S. Bak, S.B. Engelsen, Plant metabolomics: Resolution and quantification of elusive peaks in liquid chromatography-mass spectrometry profiles of complex plant extracts using multi-way decomposition methods, J Chromatogr A, 1266 (2012) 84-94.

[161] J.M. Amigo, M.J. Popielarz, R.M. Callejon, M.L. Morales, A.M. Troncoso, M.A. Petersen, T.B. Toldam-Andersen, Comprehensive analysis of chromatographic data by using PARAFAC2 and principal components analysis, J Chromatogr A, 1217 (2010) 4422-4429.

[162] Y. Xu, W. Cheung, C.L. Winder, W.B. Dunn, R. Goodacre, Metabolic profiling of meat: assessment of pork hygiene and contamination with Salmonella typhimurium, Analyst, 136 (2011) 508-514.

[163] C.M. Bishop, Pattern recognition and machine learning, springer New York2006.

[164] M. Barker, W. Rayens, Partial least squares for discrimination, J Chemometr, 17 (2003) 166-173.

[165] J. Trygg, S. Wold, Orthogonal projections to latent structures (O-PLS), J Chemometr, 16 (2002) 119-128.

[166] R. Madsen, T. Lundstedt, J. Trygg, Chemometrics in metabolomics-A review in human disease diagnosis, Anal Chim Acta, 659 (2010) 23-33.

[167] A. Kiss, C. Bordes, C. Buisson, F. Lasne, P. Lanteri, C. Cren-Olive, Data-handling strategies for metabonomic studies: example of the UHPLC-ESI/ToF urinary signature of tetrahydrocannabinol in humans, Anal Bioanal Chem, 406 (2014) 1209-1219.

[168] T. Verron, R. Sabatier, R. Joffre, Some theoretical properties of the O-PLS method, J Chemometr, 18 (2004) 62-68.

[169] A.H. Zhang, H. Sun, Y. Han, G.L. Yan, Y. Yuan, G.C. Song, X.X. Yuan, N. Xie, X.J. Wang, Ultraperformance Liquid Chromatography-Mass Spectrometry Based Comprehensive Metabolomics Combined with Pattern Recognition and Network Analysis Methods for Characterization of Metabolites and Metabolic Pathways from Biological Data Sets, Analytical Chemistry, 85 (2013) 7606-7612.

[170] B. Dieme, S. Mavel, H. Blasco, G. Tripi, F. Bonnet-Brilhault, J. Malvy, C. Bocca, C.R. Andres, L. Nada-Desbarats, P. Emond, Metabolomics Study of Urine in Autism Spectrum Disorders Using a Multiplatform Analytical Methodology, J Proteome Res, 14 (2015) 5273-5282.

[171] J. Hadrevi, M. Bjorklund, E. Kosek, S. Hallgren, H. Antti, M. Fahlstrom, F. Hellstrom, Systemic differences in serum metabolome: a cross sectional comparison of women with localised and widespread pain and controls, Scientific Reports, 5 (2015).

[172] J. Shawe-Taylor, N. Cristianini, Kernel methods for pattern analysis, Cambridge university press2004.

[173] D.S. Cao, M.M. Zeng, L.Z. Yi, B. Wang, Q.S. Xu, Q.N. Hu, L.X. Zhang, H.M. Lu, Y.Z. Liang, A novel kernel Fisher discriminant analysis: constructing informative kernel by decision tree ensemble for

1324    metabolomics data analysis, Anal Chim Acta, 706 (2011) 97-104.

1325    [174] B. Walczak, D. Massart, The radial basis functions—partial least squares approach as a flexible

1326    non-linear regression technique, Anal Chim Acta, 331 (1996) 177-185.

1327    [175] M. Bylesjo, M. Rantalainen, J. Nicholson, E. Holmes, J. Trygg, K-OPLS package: Kernel-based

1328    orthogonal projections to latent structures for prediction and interpretation in feature space, Bmc

1329    Bioinformatics, 9 (2008) 106.

1330    [176] V. Vapnik, Statistical Learning Theory, John Willey & Sons, New York, 1998.

1331    [177] H.D. Li, Y.Z. Liang, Q.S. Xu, Support vector machines and its applications in chemistry,

1332    Chemometr Intell Lab, 95 (2009) 188-198.

1333    [178] J. Luts, F. Ojeda, R. Van de Plas, B. De Moor, S. Van Huffel, J.A.K. Suykens, A tutorial on support

1334    vector machine-based methods for classification problems in chemometrics, Anal Chim Acta, 665

1335    (2010) 129-145.

1336    [179] C.J. Burges, A tutorial on support vector machines for pattern recognition, Data mining and

1337    knowledge discovery, 2 (1998) 121-167.

1338    [180] Y.B. Li, L. Ju, Z.G. Hou, H.Y. Deng, Z.Z. Zhang, L. Wang, Z. Yang, J. Yin, Y.J. Zhang, Screening,

1339    Verification, and Optimization of Biomarkers for Early Prediction of Cardiotoxicity Based on

1340    Metabolomics, J Proteome Res, 14 (2015) 2437-2445.

1341    [181] Y.B. Li, H.Y. Deng, L. Ju, X.X. Zhang, Z.Z. Zhang, Z. Yang, L. Wang, Z.G. Hou, Y.J. Zhang, Screening

1342    and validation for plasma biomarkers of nephrotoxicity based on metabolomics in male rats,

1343    Toxicology Research, 5 (2016) 259-267.

1344    [182] V.G. Uarrota, R. Moresco, B. Coelho, E.d.C. Nunes, L.A.M. Peruch, E.d.O. Neubert, M. Rocha, M.

1345    Maraschin, Metabolomics combined with chemometric tools (PCA, HCA, PLS-DA and SVM) for

1346    screening cassava (Manihot esculenta Crantz) roots during postharvest physiological deterioration,

1347    Food Chem, 161 (2014) 67-78.

1348    [183] B. Efron, Bootstrap methods: another look at the jackknife, The annals of statistics, (1979) 1-26.

1349    [184] B.F. Manly, Randomization, bootstrap and Monte Carlo methods in biology, CRC Press2006.

1350    [185] I.M. Scott, W. Lin, M. Liakata, J.E. Wood, C.P. Vermeer, D. Allaway, J.L. Ward, J. Draper, M.H. Beale,

1351    D.I. Corol, J.M. Baker, R.D. King, Merits of random forests emerge in evaluation of chemometric

1352    classifiers by external validation, Analytica Chimica Acta, 801 (2013) 22-33.

1353    [186] P.S. Gromski, H. Muhamadali, D.I. Ellis, Y. Xu, E. Correa, M.L. Turner, R. Goodacre, A tutorial

1354    review: Metabolomics and partial least squares-discriminant analysis - a marriage of convenience or a

1355    shotgun wedding, Anal Chim Acta, 879 (2015) 10-23.

1356    [187] R. Gao, J.H. Cheng, C.L. Fan, X.F. Shi, Y. Cao, B. Sun, H.G. Ding, C.J. Hu, F.T. Dong, X.Z. Yan, Serum

1357    Metabolomics to Identify the Liver Disease-Specific Biomarkers for the Progression of Hepatitis to

1358    Hepatocellular Carcinoma, Scientific Reports, 5 (2015).

1359    [188] J.H. Huang, L. Fu, B. Li, H.L. Xie, X.J. Zhang, Y.J. Chen, Y.H. Qin, Y.H. Wang, S.H. Zhang, H.Y. Huang,

1360    D.F. Liao, W. Wang, Distinguishing the serum metabolite profiles differences in breast cancer by gas

1361    chromatography mass spectrometry and random forest method, RSC Advances, 5 (2015)

1362    58952-58958.

1363    [189] Z. Lin, C.M.V. Goncalves, L. Dai, H.M. Lu, J.H. Huang, H.C. Ji, D.S. Wang, L.Z. Yi, Y.Z. Liang,

1364    Exploring metabolic syndrome serum profiling based on gas chromatography mass spectrometry and

1365    random forest models, Anal Chim Acta, 827 (2014) 22-27.

1366    [190] M. Stone, Cross-validatory choice and assessment of statistical predictions, Journal of the Royal

1367    Statistical Society. Series B (Methodological), (1974) 111-147.

1368   [191] S. Geisser, The predictive sample reuse method with applications, J Am Stat Assoc, 70 (1975)
1369   320-328.

1370   [192] J. Shao, Linear Model Selection by Cross-validation, J Am Stat Assoc, 88 (1993) 486-494.

1371   [193] D. Krstajic, L. Buturovic, D. Leahy, S. Thomas, Cross-validation pitfalls when selecting and
1372   assessing regression and classification models, J Cheminformatics, 6 (2014) 10.

1373   [194] J.A. Westerhuis, H.C.J. Hoefsloot, S. Smit, D.J. Vis, A.K. Smilde, E.J.J. van Velzen, J.P.M. van
1374   Duijnhoven, F.A. van Dorsten, Assessment of PLSDA cross validation, Metabolomics, 4 (2008) 81-89.

1375   [195] R.G. Brereton, Consequences of sample size, variable selection, and model validation and
1376   optimisation, for predicting classification ability from analytical data, Trac-Trend Anal Chem, 25 (2006)
1377   1103-1111.

1378   [196] H.D. Li, Y.Z. Liang, Q.S. Xu, D.S. Cao, Model population analysis for variable selection, J
1379   Chemometr, 24 (2010) 418-423.

1380   [197] B.C. Deng, Y.H. Yun, Y.Z. Liang, D.S. Cao, Q.S. Xu, L.Z. Yi, X. Huang, A new strategy to prevent
1381   over-fitting in partial least squares models based on model population analysis, Anal Chim Acta,
1382   10.1016/j.aca.2015.04.045 (2015).

1383   [198] R.D. Snee, Validation of regression models: methods and examples, Technometrics, 19 (1977)
1384   415-428.

1385   [199] R.W. Kennard, L.A. Stone, Computer Aided Design of Experiments, Technometrics, 11 (1969)
1386   137-148.

1387   [200] R.K.H. Galvao, M.C.U. Araujo, G.E. Jose, M.J.C. Pontes, E.C. Silva, T.C.B. Saldanha, A method for
1388   calibration and validation subset partitioning, Talanta, 67 (2005) 736-740.

1389   [201] Z.Z. Huang, Y.J. Chen, W. Hang, Y. Gao, L. Lin, D.Y. Li, J.C. Xing, X.M. Yan, Holistic metabonomic
1390   profiling of urine affords potential early diagnosis for bladder and kidney cancers, Metabolomics, 9
1391   (2013) 119-129.

1392   [202] S. Bovo, G. Mazzoni, D.G. Calo, G. Galimberti, F. Fanelli, M. Mezzullo, G. Schiavo, E. Scotti, A.
1393   Manisi, A.B. Samore, F. Bertolini, P. Trevisi, P. Bosi, S. Dall'Olio, U. Pagotto, L. Fontanesi, Deconstructing
1394   the pig sex metabolome: Targeted metabolomics in heavy pigs revealed sexual dimorphisms in plasma
1395   biomarkers and metabolic pathways, Journal of Animal Science, 93 (2015) 5681-5693.

1396   [203] J. Forshed, H. Idborg, S.P. Jacobsson, Evaluation of different techniques for data fusion of LC/MS
1397   and 1 H-NMR, Chemometrics and Intelligent Laboratory Systems, 85 (2007) 102-109.

1398   [204] T. Doeswijk, A. Smilde, J. Hageman, J. Westerhuis, F. van Eeuwijk, On the increase of predictive
1399   performance with high-level data fusion, Anal Chim Acta, 705 (2011) 41-47.

1400   [205] A. Smolinska, L. Blanchet, L. Coulier, K.A. Ampt, T. Luider, R.Q. Hintzen, S.S. Wijmenga, L.M.
1401   Buydens, Interpretation and visualization of non-linear data fusion in kernel space: study on
1402   metabolomic characterization of progression of multiple sclerosis, Plos One, 7 (2012).

1403   [206] R. Bro, H.J. Nielsen, F. Savorani, K. Kjeldahl, I.J. Christensen, N. Brünner, A.J. Lawaetz, Data fusion
1404   in metabolomic cancer diagnostics, Metabolomics, 9 (2013) 3-8.

1405   [207] L. Blanchet, A. Smolinska, A. Attali, M.P. Stoop, K.A. Ampt, H. van Aken, E. Suidgeest, T. Tuinstra,
1406   S.S. Wijmenga, T. Luider, Fusion of metabolomics and proteomics data for biomarkers discovery: case
1407   study on the experimental autoimmune encephalomyelitis, BMC bioinformatics, 12 (2011) 254.

1408   [208] A.R. Fernie, M. Stitt, On the discordance of metabolomics with proteomics and transcriptomics:
1409   coping with increasing complexity in logic, chemistry, and network interactions scientific
1410   correspondence, Plant Physiology, 158 (2012) 1139-1145.

1411   [209] S. Bocker, M.C. Letzel, Z. Liptak, A. Pervukhin, SIRIUS: decomposing isotope patterns for

1412    metabolite identification, Bioinformatics, 25 (2009) 218-224.

1413    [210] B. Zhou, J. Wang, H.W. Ressom, MetaboSearch: tool for mass-based metabolite identification

1414    using multiple databases, PLoS One, 7 (2012) e40096.

1415    [211] M. Gerlich, S. Neumann, MetFusion: integration of compound identification strategies, Journal

1416    of Mass Spectrometry, 48 (2013) 291-298.

1417    [212] F. Allen, A. Pon, M. Wilson, R. Greiner, D. Wishart, CFM-ID: a web server for annotation,

1418    spectrum prediction and metabolite identification from tandem mass spectra, Nucleic Acids Res, 42

1419    (2014) W94-99.

1420    [213] J. Draper, D.P. Enot, D. Parker, M. Beckmann, S. Snowdon, W. Lin, H. Zubair, Metabolite signal

1421    identification in accurate mass metabolomics data with MZedDB, an interactive m/z annotation tool

1422    utilising predicted ionisation behaviour 'rules', BMC Bioinformatics, 10 (2009) 227.

1423    [214] S.E. Stein, An integrated method for spectrum extraction and compound identification from gas

1424    chromatography/mass spectrometry data, J Am Soc Mass Spectr, 10 (1999) 770-781.

1425    [215] C. Steinbeck, Y.Q. Han, S. Kuhn, O. Horlacher, E. Luttmann, E. Willighagen, The Chemistry

1426    Development Kit (CDK): An open-source Java library for chemo- and bioinformatics, J Chem Inf Comp

1427    Sci, 43 (2003) 493-500.

1428    [216] M.A. Hall, Correlation-based feature selection for machine learning, The University of Waikato,

1429    1999.

1430    [217] M. Ben-Bassat, Pattern recognition and reduction of dimensionality, Handbook of Statistics, 2

1431    (1982) 773-910.

1432    [218] J. Liang, S. Yang, A. Winstanley, Invariant optimal feature selection: A distance discriminant and

1433    feature ranking based solution, Pattern Recognition, 41 (2008) 1429-1439.

1434    [219] L. Yu, H. Liu, Efficient feature selection via analysis of relevance and redundancy, The Journal of

1435    Machine Learning Research, 5 (2004) 1205-1224.

1436    [220] H.-D. Li, Y.-Z. Liang, Q.-S. Xu, D.-S. Cao, Key wavelengths screening using competitive adaptive

1437    reweighted sampling method for multivariate calibration, Ana. Chim. Acta, 648 (2009) 77-84.

1438    [221] M.D. Cao, B. Sitter, T.F. Bathen, A. Bofin, P.E. Lønning, S. Lundgren, I.S. Gribbestad, Predicting

1439    long-term survival and treatment response in breast cancer patients receiving neoadjuvant

1440    chemotherapy by MR metabolic profiling, NMR in Biomedicine, 25 (2012) 369-378.

1441    [222] E. Alba, J. Garcia-Nieto, L. Jourdan, E. Talbi, Gene selection in cancer classification using

1442    PSO/SVM and GA/SVM hybrid algorithms, Evolutionary Computation, 2007. CEC 2007. IEEE Congress

1443    on, 2007, pp. 284-290.

1444    [223] Y.-H. Yun, W.-T. Wang, M.-L. Tan, Y.-Z. Liang, H.-D. Li, D.-S. Cao, H.-M. Lu, Q.-S. Xu, A strategy that

1445    iteratively retains informative variables for selecting optimal variable subset in multivariate calibration,

1446    Anal. Chim. Acta, 807 (2014) 36-43.

1447    [224] Q. Mao, M. Bai, J.D. Xu, M. Kong, L.Y. Zhu, H. Zhu, Q. Wang, S.L. Li, Discrimination of leaves of

1448    Panax ginseng and P. quinquefolius by ultra high performance liquid chromatography

1449    quadrupole/time-of-flight mass spectrometry based metabolomics approach, Journal of

1450    Pharmaceutical and Biomedical Analysis, 97 (2014) 129-140.

1451    [225] J.S. Wang, T. Reijmers, L.J. Chen, R. Van der Heijden, M. Wang, S.Q. Peng, T. Hankemeier, G.W. Xu,

1452    J. Van der Greef, Systems toxicology study of doxorubicin on rats using ultra performance liquid

1453    chromatography coupled with mass spectrometry based metabolomics, Metabolomics, 5 (2009)

1454    407-418.

1455    [226] H.H.M. Draisma, T.H. Reijmers, J.J. Meulman, J. van der Greef, T. Hankemeier, D.I. Boomsma,

50

Hierarchical clustering analysis of blood plasma lipidomics profiles from mono- and dizygotic twin families, Eur J Hum Genet, 21 (2013) 95-101.

[227] H.-P. Kriegel, P. Kröger, A. Zimek, Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering, ACM Transactions on Knowledge Discovery from Data (TKDD), 3 (2009) 1.

[228] L. Vaclavik, A. Schreiber, O. Lacina, T. Cajka, J. Hajslova, Liquid chromatography–mass spectrometry-based metabolomics for authenticity assessment of fruit juices, Metabolomics, 8 (2012) 793-803.

[229] M.L. Ouyang, Z.M. Zhang, C. Chen, X.B. Liu, Y.Z. Liang, Application of sparse linear discriminant analysis for metabolomics data, Anal Methods-Uk, 6 (2014) 9037-9044.

[230] L.C. Phua, C.H. Wilder-Smith, Y.M. Tan, T. Gopalakrishnan, R.K. Wong, X.H. Li, M.E. Kan, J. Lu, A. Keshavarzian, E.C.Y. Chan, Gastrointestinal Symptoms and Altered Intestinal Permeability Induced by Combat Training Are Associated with Distinct Metabotypic Changes, J Proteome Res, 14 (2015) 4734-4742.

[231] T. Rajalahti, R. Arneberg, A.C. Kroksveen, M. Berle, K.M. Myhr, O.M. Kvalheim, Discriminating Variable Test and Selectivity Ratio Plot: Quantitative Tools for Interpretation and Variable (Biomarker) Selection in Complex Spectral or Chromatographic Profiles, Analytical Chemistry, 81 (2009) 2581-2590.

[232] E.C.Y. Chan, P.K. Koh, M. Mal, P.Y. Cheah, K.W. Eu, A. Backshall, R. Cavill, J.K. Nicholson, H.C. Keun, Metabolic Profiling of Human Colorectal Cancer Using High-Resolution Magic Angle Spinning Nuclear Magnetic Resonance (HR-MAS NMR) Spectroscopy and Gas Chromatography Mass Spectrometry (GC/MS), J Proteome Res, 8 (2009) 352-361.

[233] X. Lin, Q. Wang, P. Yin, L. Tang, Y. Tan, H. Li, K. Yan, G. Xu, A method for handling metabonomics data from liquid chromatography/mass spectrometry: combinational use of support vector machine recursive feature elimination, genetic algorithm and random forest for feature selection, Metabolomics, 7 (2011) 549-558.

[234] S. Mahadevan, S.L. Shah, T.J. Marrie, C.M. Slupsky, Analysis of metabolomic data using support vector machines, Analytical Chemistry, 80 (2008) 7562-7570.

[235] Y. Liu, Z. Hong, G. Tan, X. Dong, G. Yang, L. Zhao, X. Chen, Z. Zhu, Z. Lou, B. Qian, G. Zhang, Y. Chai, NMR and LC/MS-based global metabolomics to identify serum biomarkers differentiating hepatocellular carcinoma from liver cirrhosis, International Journal of Cancer, 135 (2014) 658-668.

1500 **Figure legend**

1501 **Fig.1.** Flowchart of data analysis in metabolomics.

1502 **Fig.2.** GC-MS total ion chromatograms (TICs) of tangerine peels, before (A) and after

1503 peak alignment (B). Retention time shifts in different samples were removed

1504 successfully by the multi-scale peak alignment (MSPA) approach.

1505 **Fig.3.** Concept and outline of model population analysis (MPA).

1506 **Fig.4.** Prediction error distribution of an informative, uninformative, or interfering

1507 variable before (white) and after permutation (gray) of 1000 times. Random sampling

1508 is employed. (A). Informative variable; prediction error increases after permutation.

1509 (B). Uninformative variable; prediction error shows no significant difference before

1510 and after permutation. (C). Interfering variable; prediction error may decrease after

1511 permutation.

1512

1513 Table 1. Available databases and libraries for metabolite identification

| Name | Access[a] | Current Size[b] | Website |
|---|---|---|---|
| **MS Spectral Library** | | | |
| NIST 14 | c | 276,248(242,466) | http://www.nist.gov/srd/nist1a.cfm |
| Wiley Registry of Mass Spectral Data | c | 670,000(570,000) | http://onlinelibrary.wiley.com/book/10.1002/9780470175217 |
| GolmMetabolome Database[RT] | d | 26,587 | http://gmd.mpimp-golm.mpg.de/ |
| FiehnLib | d | 1000 | http://fiehnlab.ucdavis.edu/projects/FiehnLib/index_html |
| MassBank | d | 40,889 | http://www.massbank.jp/ |
| NIST MS/MS Library | c | 234,284(9390) | http://www.nist.gov/srd/nist1a.cfm |
| ReSpect | d | 9017 | http://spectra.psc.riken.jp/ |
| METLIN | w | 71,808 | http://metlin.scripps.edu |
| **Chemical Substance Database** | | | |
| PubChem Compound Dabatase | d | >53 million | http://www.ncbi.nlm.nih.gov/pccompound |
| ChemSpider | w | >21 million | http://www.chemspider.com/ |
| Manchester Metabolomics Database | d | 42,553 | http://dbkgroup.org/MMD/ |
| BiGG Database | w | 2835 | http://bigg.ucsd.edu/bigg |
| BioCyc (MetaCyc) | | UNKNOWN | http://biocyc.org/ |
| CAS Registry | c | >89 million | http://www.cas.org/ |
| CSLS | w | UNKNOWN | http://cactus.nci.nih.gov/ |
| GDB databases | d | ~166 billion | http://www.gdb.unibe.ch/gdb/ |
| Dictionary of Natural | c | 240,007 | http://dnp.chemnetbase.com/dictionary-search.do?method=view&id=1079994 |

| Products | | | 5&struct=start&props=&&si= |
|---|---|---|---|
| Beilstein database | c | >500 million | http://www.elsevier.com/online-tools/reaxys |
| KEGG ligand database | d | 17,282 | http://www.genome.jp/kegg/ligand.html |
| ChEBI | d | 40,211 | http://www.ebi.ac.uk/chebi/ |
| HMDB | d | 41,806 | http://www.hmdb.ca/ |
| KNApSAcK | d | 50,899 | http://kanaya.naist.jp/KNApSAcK/ |
| LIPID MAPS | d | 37,566 | http://www.lipidmaps.org/ |
| LipidBank | w | 7,009 | http://www.lipidbank.jp/ |
| METLIN | w | 240,501 | http://metlin.scripps.edu |
| SDBS | w | 34,000 | http://sdbs.db.aist.go.jp/sdbs/cgi-bin/cre_index.cgi |

1514 [a]Access right to the database, c, d and w denote commercial, downloadable and online access,

1515 respectively.

1516 [b] Number of unique compounds for corresponding library are provided in the bracket.

1517 [RT] Retention indices are included.

1518

1519 Table 2. Available metabolite identification tools and related tools assisting metabolite

1520 identification

| Name | Reference | Website |
|---|---|---|
| **GC-MS Spectrum Identification** | | |
| MassLib | | http://www.masslib.com/[c] |
| MOLGEN-MS | [79] | http://molgen.de/?src=documents/molgenms.html[d,w] |
| Mass Spectrum Interpreter | [81] | http://chemdata.nist.gov/mass-spc/interpreter/[d] |
| **Accurate Mass** | | |
| MetWorks | | http://www.thermoscientific.com[c] |
| MetabolitePilot | | http://www.absciex.com[c] |
| Seven Golden Rules | [88] | http://fiehnlab.ucdavis.edu/projects/Seven_Golden_Rules/[d] |
| SIRIUS | [209] | http://bio.informatik.uni-jena.de/sirius2/[d] |
| MI-Pack | [103] | http://www.biosciences-labs.bham.ac.uk/viant/mipack[d] |
| MetaboSearch | [210] | http://omics.georgetown.edu/metabosearch.html[d] |
| **MS/MS Spectrum Prediction** | | |
| Mass Frontier | | http://www.thermoscientific.com[c,g] |
| ACD/MS Fragmenter | | http://www.acdlabs.com[c,g] |
| MetISIS | [115] | http://omics.pnl.gov/software/[d] |
| MyCompoundID | [116] | www.mycompoundid.org[w] |
| *In silico* **Fragmentation** | | |
| FiD | [118] | http://www.cs.helsinki.fi/group/sysfys/software/fragid/[d] |
| Mass-MetaSite | [119] | http://www.moldiscovery.com/software/massmetasite[c] |
| MetFrag | [120] | http://c-ruttkies.github.io/MetFrag/[d,w] |
| FingerID | [121] | https://github.com/icdishb/fingerid[d] |
| MetFusion | [211] | http://msbi.ipb-halle.de/MetFusion/[w] |

| | | |
|---|---|---|
| CFM-ID | [122, 212] | http://cfmid.wishartlab.com/[d,w] |
| *De Novo* **Analysis** | | |
| SIRIUS[2] | [123, 124] | http://bio.informatik.uni-jena.de/sirius2/[d] |
| MAGMa | [177] | http://www.emetabolomics.org/ |
| **Molecule Ion Annotation** | | |
| PUTMEDID-LCMS | [108] | http://www.mcisb.org/resources/putmedid.html[d] |
| CAMERA | [109] | http://metlin.scripps.edu/xcms/useful_links.php[d] |
| IDEOM | [110] | http://mzmatch.sourceforge.net/ideom.php[d] |
| MZedDB | [213] | http://maltese.dbs.aber.ac.uk:8888/hrmet/index.html[w] |
| MAIT | | |
| **Mass Spectra Deconvolution** | | |
| AMDIS | [214] | http://chemdata.nist.gov/dokuwiki/doku.php?id=chemdata:amdis[d] |
| DeconvolutionReporting Software | | http://www.chem.agilent.com/en-US/products-services/Software-Informatics/Deconvolution-Reporting-Software-%28DRS%29/Pages/default.aspx[c] |
| AnalyzerPro | | http://www.spectralworks.com/analyzerpro.html[c] |
| ChromaTOF® | | http://www.leco.com/products/separation-science/software-accessories/chromatof-software[c] |
| **Formula Generation** | | |
| OMG | [77] | http://sourceforge.net/projects/openmg/[d] |
| The Chemistry Development Kit | [215] | http://sourceforge.net/projects/cdk/[d] |
| Formula To Mass To Formula | | http://www.ch.ic.ac.uk/java/applets/f2m2f/[w] |
| Molecular Formula finder | | http://www.chemcalc.org/mf_finder[w] |
| HiRes | | http://hires.sourceforge.net/[w,d] |

1521 [c]Commercially available. [d]Freely downloadable to the local site. [w]Freely accessed via web

1522 interface. [g]Also suitable for GC-MS spectrum.

1523

1524

Table 3.　A taxonomy of variable selection techniques with the mentioned methods

| Methods | Classifier | Interpretability | Consider the interaction effect among variables or not | Variable ranking or subset selection | Computation speedy | Reference |
|---------|-----------|------------------|---------------------|---------------------|-----------|-----------|
| PLS-weights | PLS | Based on loading weight matrices of PLS modeling | NO | Ranking | High | [131] |
| PLS-VIP | PLS | Accumulate the importance of each variable being reflected by loading weights from each latent variable of PLS | NO | Ranking | High | [132] |
| PLS-regression coefficient | PLS | A single measure of association between each variable and the response. | NO | Ranking | High | [133] |
| Correlation | No classifier | Calculate simply between variables and classification label. | NO | Ranking | High | [216] |
| Information gain | No classifier | | NO | Ranking | High | [217] |
| Euclidean distance | No classifier | | NO | Ranking | High | [218] |
| Mutual information | No classifier | | NO | Ranking | High | [219] |
| CARS | PLS | Realize a competitive feature selection based on the absolute regression coefficients. | NO | Subset selection | High | [220] |
| GA-PLS-DA | PLS-DA | GA is used as an optimal algorithm to find the optimal subset with PLS-DA classifier. | NO | Subset selection | Low | [221] |
| PSO-SVM | SVM | PSO is used as an optimal algorithm to find the optimal subset with SVM classification method. | NO | Subset selection | Medium | [222] |
| Random Forest | Decision Tree | Rank the variables by the percent increase of misclassification error when the | YES | Ranking | Medium | [143] |

55

| | | variable is permuted randomly. | | | | |
|---|---|---|---|---|---|---|
| SPA | PLS-DA | Identify and rank the informative variable based on the difference between the prediction errors of normal and permutated subwindow for each variable. | YES | Ranking | Medium | [145] |
| MIA | SVM | Give a measure based on the difference between the prediction errors of inclusion and exclusion for each variable with the margin of SVM | YES | Ranking | Medium | [146] |
| INTERACT | No classifier | Based on inconsistency and symmetrical uncertainty measurements for finding interacting features | YES | Subset selection | High | [142] |
| VCN | PLS-DA | Compute the complementary information between variables and then effectively discover biomarker with the help of mutual associations of metabolites. | YES | Ranking | Medium | [149] |
| IRIV | PLS | Find the optimal subset of variables through observing the difference between the prediction errors of inclusion and exclusion for each variable. | YES | Subset selection | Medium | [223] |
| VISSA | PLS | Search for the optimal variable combinations through shrinking the variable space smoothly | YES | Subset selection | Medium | [148] |

1525
1526
1527

1528

Table 4.   An overview of multivariate analysis methods for modeling

| Method | Category | Advantage | Disadvantage | Applications in metabolomics |
|---|---|---|---|---|
| PCA | unsupervised | Suit to provide an overview of a large dataset. | Class information is not considered. | [224, 225] |
| HCA | unsupervised | Suit to provide an overview of the clusters of samples. | Class information is not considered. Variable importance is not obtained. | [226, 227] |
| SOM | unsupervised | Account for non-linear in the data | Class information is not considered. | [155-157] |
| PARAFAC2 | unsupervised | Can handle shifted data with baseline | Can be more sensitive to noise | [160, 162] |
| LDA | supervised | Easy and fast. Suit to linear and low dimensional data. | Not suit to high dimensional data | [228, 229] |
| PLS-DA | supervised | Particularly suit to linear and co-linear data. | Not suit to unbalanced data. | [4, 230, 231] |
| OPLS-DA | supervised | Particularly suit to linear and co-linear data. Good visualization ability and interpretation ability. | Not suit to unbalanced data. | [169-171, 232] |
| SVM | supervised | Suit to linear and nonlinear problem. High flexibility in modeling non-linear data. | Lack of transparency of the results. Model tuning is complex | [180-182, 233, 234] |
| RF | supervised | Suit to linear and nonlinear problem. Resistance to outliers. | Relatively low computation speed | [187-189, 235] |

1529

1. Pre-processing

- Noise filtering and baseline correction
- Peak detection and deconvolution
- Alignment
- Normalization

Objects: (e.g. Plants)

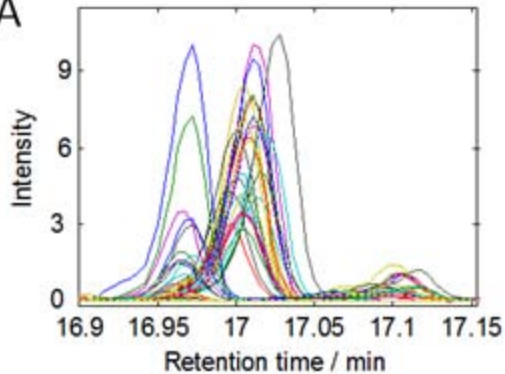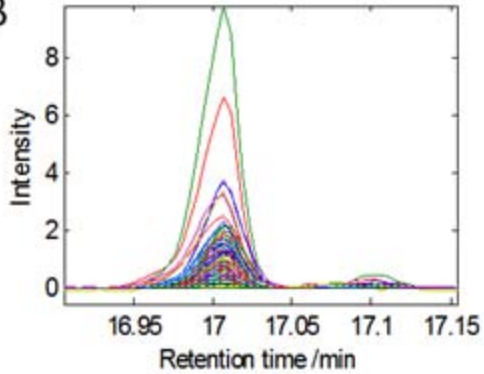Features: Measured values

4. Modeling of the data

PCA

PLS-DA

RF

SVM

→ ANN

→ Others

2. Identification of metabolites

3. Variable selection

LW

VIP

VISSA

MPA

→ IRIV

→ Others

Baichuan Deng received his M.S. and Ph.D. in chemometrics from the Department of Chemistry at the University of Bergen (Norway), in 2012 and 2015, respectively. He is currently a Distinguished Associate Professor in South China Agriculture University. His main research interests include nutritional metabolomics, chemometrics and bioinformatics.

Naiping Dong obtained both a B.S. and a Ph.D. from Central South University (China). He is now research associate in Department of Applied Biology and Chemical Technology, the Hong Kong Polytechnic University. In 2013, He obtained his Ph.D. in analysis of high throughput tandem mass spectrometry in proteomics with Prof. Yizeng Liang as supervisor. Now, his research mainly focuses on applying chemometric and statistical methods in processing chromatographic and mass spectrometric data generated from biological samples.

Shao Liu is the Professor of Xiangya Hospital, Central South University. He has published more than 100 research papers on leading scientific journals. Major research interests include: （1）Natural medicine resources and drug discovery；（2）Quality control of traditional Chinese medicine and metabolomics based on chemometrics.

Yizeng Liang is a professor of chemometrics and analytical chemistry of College of Chemistry and Chemical Engineering, Central South University, China. Editor of "Chemometrics and Intelligent Laboratory Systems" (since 2007). Research interests include analytical chemistry and chemometrics; quality control of traditional Chinese medicines; metabolomics and proteomics; Data mining in chemistry and Chinese medicines. Professor Liang has published more than 420 scientific research papers since 1989 in SCI source journals. Besides, he has published 10 books (8 in Chinese and 2 in English).

Dabing Ren obtained both a B.S. and a Ph.D. from Central South University (China). In 2012, he started his Ph.D. in Professor Yizeng Liang's group with a focus on countercurrent chromatography (CCC). He mainly used thermodynamic models to study the partition behavior of solutes involved in the CCC separation process. And he developed some useful methods used to correlate and predict the solute partition coefficient in biphasic solvent systems. In 2015, he started career in Kunming University of Science and Technology (China). Now, his research interest mainly covers metabolomics, chromatographic analysis and mass spectroscopy.

Lunzhao Yi is a professor of analytical chemistry and food science of Yunnan Food Safety Research Institute, Kunming University of Science and Technology, China. In 2004, she started her Ph.D. in Professor Yizeng Liang's group with a focus on chemometrics and metabolomics. In 2007, she started career in Central South University (China) until 2014. Her research interests include analytical chemistry, chemometrics, metabolomics and food chemistry. Professor Yi has published more than 50 scientific research papers since 2006 in SCI source journals.

Yonghuan Yun received his B.S. in Pharmaceutical Engineering from Central South University (China) in 2011. He is currently pursuing his Ph.D. in Analytical Chemistry at Central South University under the supervision of Professor Yizeng Liang. His research is focused on developing new algorithm of chemometrics and bioinformatics in the field of near infrared and Raman spectroscopy, metabolomics and genomics.