

Accepted Manuscript

This document is the Accepted Manuscript version of a Published Work that appeared in final form in *Analytical Chemistry*, copyright © American Chemical Society after peer review and technical editing by the publisher.

To access the final edited and published work see
<http://dx.doi.org/10.1021/acs.analchem.7b00743>

Saer Samanipour, Malcolm J. Reid, Kevin V. Thomas. 2017.
Statistical variable selection: an alternative prioritization
strategy during the nontarget analysis of LC-HR-MS data.
Analytical Chemistry . 89(10): 5585-5591. ISSN 0003-2700.

It is recommended to use the published version for citation.

Statistical variable selection: An alternative prioritization strategy during the non-target analysis of LC-HR-MS data

Saer Samanipour,^{*,†} Malcolm J. Reid,[†] and Kevin V. Thomas^{†,‡}

[†]*Norwegian Institute for Water Research (NIVA), 0349 Oslo, Norway*

[‡]*Queensland Alliance for Environmental Health Science (QAEHS), University of Queensland, 39 Kessels Road, Coopers Plains QLD 4108, Australia*

E-mail: saer.samanipour@niva.no

Abstract

Liquid chromatography coupled to high resolution mass spectrometry (LC-HR-MS) has been one of the main analytical tools for the analysis of small polar organic pollutants in the environment. LC-HR-MS typically produces a large amount of data for a single chromatogram. The analyst is therefore required to perform prioritization prior to non-target structural elucidation. In the present study we have combined the F-ratio statistical variable selection and the apex detection algorithms in order to perform prioritization in data sets produced via LC-HR-MS. The approach was validated through the use of semi-synthetic data, which was a combination of real environmental data and the artificially added signal of 31 alkanes in that sample. We evaluated the performance of this method as a function of four false detection probabilities namely: 0.01, 0.02, 0.05, and 0.1%. We generated 100 different semi-synthetic data sets for each F-ratio and evaluated that data set using this method. This design of experiment created a population of 30,000 true positives and 32,000

15 true negatives for each F-ratio, which was considered sufficiently large enough in order
16 to fully validate this method for analysis of LC-HR-MS data. The effect of both the F-
17 ratio and signal to noise ratio (S/N) on the performance of the suggested approach were
18 evaluated through normalized statistical tests. We also compared this method to the
19 pixel-by-pixel as well as peak list approaches. More than 92% of features present in the
20 final feature list via F-ratio method were also present in conventional peak list generated
21 by MZmine. However, this method was the only approach successful in classification of
22 samples, thus prioritization, when compared to the other evaluated approaches. The
23 application potential and limitations of the suggested method discussed.

24 Introduction

25 A large number of small polar organic pollutants are considered as chemicals of emerging
26 concern (CECs) due to their fate and behavior in the environment (as reviewed by Klečka
27 et al.¹ and La Farre et al.²). Liquid chromatography coupled to the high resolution mass
28 spectrometry (LC-HR-MS) has become the leading analytical instrumentation for analysis
29 of these pollutants in different environmental compartments.³⁻⁵ Measuring these pollutants
30 in the environment takes place through three different and/or complementary approaches,
31 namely target analysis, suspect analysis, and non-target analysis.^{4,6-8} For target analysis
32 the analyst has all the necessary information, including the retention time and the spec-
33 tral information, for confident identification of a target analyte in complex environmental
34 samples.⁹ On the other hand for the suspect screening only limited information is available
35 while during non-target analysis the analyst does not have any prior information regarding
36 the identity of the analytes in the sample.^{7,9} Even though non-target analysis is the most
37 difficult and the least certain of the three mentioned approaches, this method is essential for
38 the discovery of new CECs in the environment.^{4,6,7}

39
40 Confident identification of pollutants based only on the data generated via non-target

41 analysis on LC-HR-MS, is a challenging task due to the volume and complexity of the data
42 .^{10,11} During the non-target analysis, each sample may produce thousands of features, where
43 each includes a measured exact mass, intensity, and the retention time.^{12,13} Therefore, the
44 analyst may have to prioritize among the features for structural elucidation. For LC-HR-
45 MS data, there have been different approaches used for prioritization during the non-target
46 analysis. The simplest approach applies the absolute intensity and the detection frequency
47 as the main criteria for prioritization.^{9,12} However, high signal intensity and the detection
48 frequency does not guaranty environmental relevance. Another approach utilizes either tox-
49 icity information (through effect-directed analysis¹⁴) or the elemental composition (through,
50 for example, filters for halogenated compounds¹²). However, the mentioned approaches are
51 complicated and may be biased towards a certain family of compounds, for example halo-
52 genated ones. A less used approach, particularly in the field of environmental analysis, has
53 been the application of unsupervised and/or supervised statistical methods, such as princi-
54 pal component analysis (PCA) and partial least square discrimination analysis (PLS-DA)
55 for prioritization of the relevant features.^{13,15} These statistical methods perform well when
56 used in metabolomics due to a more clear change in the sample composition. However, these
57 same methods may suffer when there is a high level of redundancy/similarity in the analyzed
58 samples.^{16,17} Recent studies have shown the superior performance of the supervised F-ratio
59 method combined with PCA for analysis of the data via gas chromatography coupled to
60 low resolution mass spectrometry (GC-MS) of complex samples.¹⁶⁻²⁰ However, the F-ratio
61 method has never been used/optimized for the non-target analysis of the data recorded via
62 LC-HR-MS and particularly for complex environmental samples.

63

64 The mentioned statistical variable selection approaches can be applied to either a peak
65 list^{11,13,15} or the whole chromatogram.¹⁶⁻²⁴ Even though processing of the peak list is faster
66 than the whole chromatogram due to its smaller size compared to the whole chromatogram,
67 the raw chromatogram must go through preprocessing steps such as signal deconvolution,

68 peak finding, peak picking, and peak integration in order to generate a final peak list useful
69 to a prioritization method. All these preprocessing steps are prone to error when dealing with
70 highly complex samples.²²⁻²⁵ The application of the statistical variable selection approaches
71 to the whole chromatogram has been shown to result in reliable models and therefore, reli-
72 able prioritization.¹⁶⁻²⁵

73

74 The aim of this study is to adapt, comprehensively validate, and test the applicability
75 of the F-ratio method for the non-target analysis of LC-HR-MS chromatograms of complex
76 environmental samples. The F-ratio was applied to the whole chromatogram in order to
77 minimize the data manipulation and produce a reliable statistical model. We combined the
78 F-ratio method with the apex detection as well as adduct and isotope removal algorithms,
79 in order to adapt this method to be used for non-target analysis of LC-HR-MS data. We
80 comprehensively validated this method using a semi-synthetic data set, which consisted of
81 the background signal generated from the real environmental samples with the addition of
82 the signal of 31 alkanes randomly distributed as true positives and true negatives, and noise.
83 This data set was evaluated 400 times where the random selection of the alkanes and the
84 background signal caused generation of a completely different sample for each evaluation.
85 Finally, the chromatograms of 15 sludge extracts from three different locations in Norway
86 and three blanks were analyzed using the F-ratio method as well as conventional peak picking
87 algorithms. We also applied the F-ratio method to the peak list and compared this feature
88 list to the one produced via using the whole chromatogram. The feature lists via F-ratio
89 were compared to the peak lists generated by a conventional peak pick method, in order to
90 further evaluate and/or validate the applicability of this method for non-target analysis of
91 LC-HR-MS data.

92 **Experimental Methods**

93 **Environmental Sampling and Sample Preparation**

94 15 sludge samples from three different wastewater treatment plants (WWTP) in Norway were
95 collected (i.e. five replicates for each WWTP), during the spring of 2015. These WWTPs
96 were located at Oslo, Hamar, and Gjøvik. The plants in Oslo and Hamar were equipped
97 with a three stage treatment process, including physical, chemical, and biological treatment
98 whereas the plant in Gjøvik has only the physical and chemical treatments. More details on
99 the chemicals, the suppliers, and the sample preparation steps are provided in section S1 of
100 Supporting Information.

101 **Instrumental Conditions and Analysis**

102 All the extracts were analyzed employing, Waters Acquity UPLC system (Waters Milford,
103 MA, USA). An Acquity UPLC HSS C18 column (2.1×150 mm, particle size 1.8 μ m) (Wa-
104 ters, Milford, MA, USA) was used for all the separations. A mixture of solvent A, 5 mM
105 ammonium formate at pH 3.0 and solvent B, acetonitrile with 0.1% formic acid at a constant
106 flow rate of 0.4 ml min^{-1} was used for the chromatographic separations. The gradient varied
107 from 87% of solvent A to 5% of solvent B. More details regarding this method are provided
108 elsewhere.²⁶ Both the analytical column and the column guard were kept as $50 \text{ }^\circ\text{C}$ during
109 the separations.

110
111 Xevo G2-S Q-TOF-MS (Waters Milford, MA, US) was used for analysis of all 18 samples,
112 including the 15 sludge extracts and 3 blanks. The MS¹, with a collision energy of 6 eV, and
113 the MS², with a collision energy ramp between 15 to 45 eV were simultaneously recorded
114 during the whole chromatogram. We employed a mass range of between 90 Da and 700 Da
115 with a sampling frequency of 1.8 Hz. More information regarding the mass spectrometer
116 conditions is available elsewhere.²⁶

117 Data Processing and Workflow

118 All the chromatograms were recorded in profile mode employing MassLynx (Waters Milford,
119 MA, US). The chromatograms were then exported as netCDF files, using DataBridge package
120 (Waters Milford, MA, US) incorporated in the MassLynx software. The MS¹ channel (i.e.
121 low collision energy channel 6 eV) chromatograms were used for the data analysis. All the
122 exported chromatograms were then imported into Matlab.²⁷ These files were processed using
123 the following sequence of steps in the stated order including the binning, retention alignment,
124 data matrix generation, F-ratio calculation, null-distribution validation, zero mask applica-
125 tion, chromatogram folding, apex detection, and finally adducts and isotope removal. All the
126 steps taken in this workflow are explained in detail below and section S2 of the Supporting
127 Information.

128

129 We binned the exported chromatograms using a bin thickness of 10 mDa, which was
130 based on the observed mass accuracy of ± 5 mDa in our data set (section S2.1). The mass
131 accuracy was defined based on the shift observed in the measured mass of the calibrant
132 injected every 20 s into the source. The binned chromatograms were then retention aligned
133 with a home-developed algorithm inspired by the piecewise method previously developed and
134 validated by Synovec group.^{25,28} We added an additional mass spectral correlation control
135 in order to increase the accuracy of the retention alignment. More details regarding both
136 binning and retention alignment processes are provided in Supporting Information section
137 S2. The retention aligned chromatograms were then unfolded to create a long vector of in-
138 tensities for every single measured m/z value. These vectors were then stacked on top of each
139 other in order to produce a large matrix which was used for the statistical prioritization.
140 Every row in this matrix was a sample while every column was an independent variable.
141 The F-ratio was calculated for each variable,¹⁶ or column of the matrix, based on *a priori*
142 knowledge of the sample classification (section S2.3). An F-ratio threshold was calculated
143 using the probability distribution generated via null-distribution analysis.¹⁹ This procedure

144 aims to minimize the number of false positive detections as well as the method validation
145 during the analysis (see section S2.4). The variables that had an F-ratio smaller than the
146 defined threshold, based on the null-distribution, were set to zero in the data matrix. This
147 process was referred to as the zero mask application. Each zero mask applied chromatogram
148 then was folded back into a matrix where a row was one scan and a column was the signal
149 for a m/z value. We performed apex detection in the folded chromatograms (see section
150 S2.5 of Supporting Information). The apex detection groups the non-zero and statistically
151 meaningful variables which can be represented as a feature in the chromatogram. For ex-
152 ample all the non-zero variables in a chromatographic peak can be grouped and represented
153 via only one pair of retention time and m/z value, thus a feature. Therefore, the apex de-
154 tection generates a list of unique retention time and m/z value pairs for each sample. This
155 differs from conventional peak picking algorithms in that apex detection does not perform
156 signal modeling and/or integration therefore minimizes the signal manipulation. Finally,
157 the adducts and the isotopes were removed from this list in order to create the final unique
158 feature list for each chromatogram. This workflow provides the necessary initial information
159 for discovery-based non-target analysis of complex samples analyzed via LC-HR-MS.

160

161 We also performed F-ratio analysis on the peak list produced by conventional peak pick-
162 ing algorithm, MZmine 2²⁹(explained in detail below). The peak list was retention-aligned
163 using a home-developed method using a mass window of 2 mDa and a retention window of 2
164 S. The retention aligned peak tables were used for F-ratio and null-distribution calculations.
165 The peaks in the peak list with an F-ratio larger than the threshold were kept in order
166 to produce the feature list. The feature list, finally, was processed for adduct and isotope
167 removal in order to generate the final feature list.

168

169 **Data pretreatment**

170 During the validation process of the F-ratio method (i.e. analysis of the semi-synthetic data),
171 we did not employ any data pre-treatment methods such as mean-centering, standardiza-
172 tion, and normalization. This choice enabled us to comprehensively evaluate the effect of
173 introduced noise on the performance of the F-ratio method. For the environmental sample
174 analysis, we tested different data pre-treatment methods such as mean-centering, standard-
175 ization, and normalization before processing the data set with F-ratio method. However,
176 these pre-treatments did not affect the final unique feature list for the analyzed data set.
177 Therefore, we decided to work with the raw data and avoid performing any type of pre-
178 treatment.

179 **Computations**

180 All the mentioned data processing steps were performed via Matlab, employing a Windows
181 7 Professional version (Microsoft Inc, USA) workstation computer with 12 CPUs and 128
182 GB of memory.

183 **MZmine Peak Picking**

184 The conventional peak list for each chromatogram was generated using MZmine 2.²⁹ The
185 peak picking was performed by mass detection followed by GridMass 2D peak detection. A
186 five scan window was selected for the smoothing of the chromatogram in the time dimension
187 and a 10 mD window was used in the mass dimension. A minimum signal of 300 counts was
188 required for a peak to be considered as a meaningful peak. These parameters were optimized
189 based on the observed mass accuracy and the peak widths in both time and mass domains.
190 These parameter settings resulted in feature numbers varying between 7,500 for blanks and
191 12,500 for the samples from Oslo WWTP.

192 **Principal Component Analysis (PCA)**

193 We employed principal component analysis (PCA)³⁰ for classification/separation of the sam-
194 ple groups. We performed PCA on the peak list generated via MZmine, variable selected peak
195 list (i.e. the F-ratio applied to the peak list generated via MZmine) with F-ratio method,
196 the whole chromatogram (i.e. pixel-by-pixel), and the variable selected chromatogram em-
197 ploying F-ratio method. The PCA was performed on the mean centered data utilizing the
198 singular value decomposition algorithm.³⁰

199 **Results and Discussion**

200 We validated the F-ratio method for data generated via LC-HR-MS, employing both semi-
201 synthetic data and the real environmental data. The use of semi-synthetic data enabled us
202 to perform a large number of evaluations (i.e. total number of detection cases 62,000×4)
203 knowing exactly the added signal, noise, and relative intensity of the added signal which,
204 translated into comprehensive validation of the proposed method. This would not have
205 been possible using spiked samples due to the limitation on the number of standards and
206 injections as well as the potential interference between the sample and the standard mixture.
207 This study is the first implication of this method for the data generated via LC-HR-MS as
208 well as adaptation of this method in order for it to be included in non-target identification
209 workflows. This method enables the direct prioritization of the unique features, which are
210 the main cause of the separation of different sample groups. Therefore, the identification
211 efforts can be focused on the prioritized unique features.

212 **Validation via Semi-synthetic Data**

213 We employed a semi-synthetic data set, which consisted of a combination of real environ-
214 mental data and synthetic signal, for comprehensive validation of the F-ratio method. The
215 signal of 31 alkanes (i.e. the neutral monoisotopic masses, Table S1) was added at different

216 concentrations to a background signal, which came from the real environmental data. During
217 each analysis, these 31 alkanes were divided in two randomly selected groups where the first
218 group of 15 alkanes was added to the background signal at concentration levels that were
219 statistically meaningful. Therefore, for these 15 alkanes the resulting F-ratios were larger
220 than the threshold. For the second group, 16 alkanes were added to the background at a sta-
221 tistically constant concentration. Four different F-ratios of 208, 30, 28, and 13 having false
222 positive detection probabilities of 0.01, 0.02, 0.05, and 0.1 %, respectively, were evaluated.
223 Each F-ratio value was evaluated 100 times with different: background signal, combination
224 of alkanes, concentration levels, and retention times of true positives (i.e. 15 alkanes) and
225 true negatives, thus a total of 400 evaluations. The generation of the these semi-synthetic
226 data is described in detail in Supporting Information, section S3. Alkanes were selected for
227 our analysis because these compounds are not ionized by ESI source therefore we were sure
228 that these compounds were not present in the real background signal. This design of exper-
229 iment created a set of 15 true positives, 16 true negatives, and a different background signal
230 during each evaluation, which enabled us to comprehensively examine the capabilities and
231 limitations of the F-ratio method. The number of evaluation (i.e. 100 for each F-ratio) was
232 selected based on our preliminary assessment, that showed that 100 analysis for each F-ratio
233 would generate a large enough population of true positives (TP) 30,000 and true negatives
234 (TN) 32,000 for that F-ratio, in order to fully validate this method. To compare the effect
235 of different F-ratio probability value on the final results, we employed normalized statistical
236 parameters such as rate of false positive, rate of false negative, sensitivity, specificity, and
237 accuracy.³¹

238

239 Increased F-ratios resulted in a smaller number of false positives and a larger number of
240 false negatives. The number of false positive detection ranged from 2,518 cases for the F-ratio
241 of 208 having a probability of 0.01% to 9,525 cases for the F-ratio of 13 with a probability
242 of 0.1%, Table 1 and Figure S7. The largest number of false negative detections of 2204 was

243 observed for an F-ratio of 208 whereas the smallest number of false negative detections of
244 1,404 was caused by an F-ratio of 13, Table 1. These trends were due to the fact that the
245 selection of a large F-ratio value (i.e. more strict selection criterion) lowers the probability
246 of false positive detection while increasing the probability of false negative detections. The
247 observed changes in F-ratio method performance were better projected through normalized
248 statistical parameters such as rate of false positive detection, rate of false negative detection,
249 sensitivity, specificity, and accuracy, Table 1. For example, the drop in the specificity and
250 accuracy observed for F-ratio of 13 (probability of 0.1%) showed the inadequacy of this F-
251 ratio for the analyzed data set (Figure S7). This drop also indicated that this F-ratio may
252 cause a large number of false positive detections when analyzing this data set. Therefore,
253 the analyst is required to find an optimized F-ratio value in order to minimize the number
254 of potential false positive detection while limiting the number of false negatives. Among the
255 four F-ratios evaluated, the value of 28 (probability of 0.05%) showed to be the optimized
256 one, considering that this value provided the largest accuracy level, second largest sensitivity
257 level while maintaining a high level of specificity (Figure S7).

258

259 Further evaluation of our data set, showed that for F-ratios ≥ 28 more than 70% of the
260 false positive cases were coming from the background signal rather than the true negatives
261 (i.e. added signal of alkanes at constant concentrations). The observed trend was caused
262 by the high level of variability artificially introduced into the background signal during the
263 background generation. For F-ratio of 13, around 50% of the false positives were true neg-
264 atives. In this case even though the signal of true negatives did not have a large level of
265 variability between sample groups, once added to the background, the variability in that
266 signal increased due to the inherent large variance in the background signal. Therefore,
267 these true negatives produced a large enough F-ratio, which met the F-ratio threshold and
268 were selected as positive detections. When looking at the false negative cases, for all four
269 F-ratios, the main causes of false negative detection were the large variability in the back-

270 ground signal and the S/N threshold setting during the apex detection. Also in this case,
271 the random variability introduced into the true positives signal was not large enough to
272 overcome the variability present in the background signal. We tested these hypothesis by
273 increasing the initial concentration of the added signal of the true positives from 5% to 15%
274 and also increasing the concentration factor from 2-8 to 2-20 (see section S3 of SI for more
275 details). With an F-ratio of 28, increasing these parameters reduced drastically the number
276 of false positive detection from 2,864 to 253 cases as well as the number of false negatives
277 from 1,570 to 35 cases after 100 simulations. These results indicated that the combination
278 of low level concentration of added alkanes, their low between sample group variability, and
279 finally the large level of variability introduced into the background signal have an important
280 effect on the performance of this algorithm.

281

282 Mean centering and standardization (i.e. division by the square root of standard devia-
283 tion of each variable) with an F-ratio of 28, added signal of 5%, and the concentration factor
284 of 2-8 reduced the number of both false positive detection and false negative detection from
285 2,864 to 350 cases and from 1,570 to 97 cases, respectively after 100 simulations. These
286 pre-treatments' approaches both decrease the noise levels in the data set while emphasizing
287 the underlying trend.³⁰ This implies that these data pre-treatments reduced the effect of
288 artificially introduced noise in the data set while emphasizing the between group variability,
289 thus a decrease in the number of false positive and false negative detection. The type of data
290 pre-treatments employed prior to the F-ratio analysis is data set and objective dependent.³⁰
291 Therefore, the analyst is required to optimize these data pre-treatments approaches in ad-
292 vance in order to be able to produce reliable results. Further investigation on the effect of
293 these parameters on the F-ratio method are needed and will be subject of our future studies.

Table 1: The number of false positive detections, number of false negative detections, rate of false positive, rate of false negative, sensitivity, specificity, and accuracy parameters calculated for four different F-ratio values based on 100 evaluations for each F-ratio probability value.

Parameter	F-ratio values (probability of false positive detection %)			
	208 (0.01)	30 (0.02)	28 (0.05)	13 (0.10)
False positive detection ^a (FP)	2,518	3,172	2,864	9,525
False negative detection ^a (FN)	2,204	2,220	1,570	1,404
Rate of false positive detection ^b (%)	7.3	9.0	8.0	23.0
Rate of false negative detection ^c (%)	6.8	6.9	5.0	4.5
Sensitivity ^d (%)	93.2	93.1	95.0	95.5
Specificity ^e (%)	92.7	91.0	91.8	77.1
Accuracy ^f (%)	92.9	92.0	93.3	85.0

^aThis parameter represents the number false positive detection out of total number of detections of 62,000, including 30,000 true positives (TP) and 32,000 true negatives (TN); ^bThe rate of false positive³¹ was calculated as $FP/(FP+TN)$; ^cThe rate of false negative³¹ was calculated as: $FN/(TP+FN)$; ^dThe sensitivity³¹ values were calculated using: $TP/(TP+FN)$; ^eThe specificity³¹ values were calculated with: $TN/(TN+FP)$; ^fThe accuracy³¹ values were calculated employing: $(TP+TN)/(TP+FP+FN+TN)$.

294 The effect of S/N on F-ratio algorithm

295 The S/N is an important parameter, which affects the performance of the F-ratio algorithm
 296 particularly during the apex detection. The apex detection step aims to reduce the level
 297 of redundancy in the data set by grouping variables, that can be represented by a unique
 298 one (see section S2.5 for more detailed information). We evaluated the effect of S/N on the
 299 results of the algorithm with an F-ratio of 28. This evaluation was performed by varying
 300 the S/N from 1 to 10 (i.e. 1, 3, and 10) and performing 20 analysis for each S/N value. The
 301 F-ratio of 28 was selected based on the fact that it appeared to be the optimized F-ratio for
 302 the evaluated data set.

303

304 We observed a slight decrease in the number of false positive detection as a function of
 305 increase in the S/N while the increase in the S/N had a positive effect on the number of
 306 detected false negatives, Table 2. However, the changes in the S/N did not appear to cause
 307 a large variation in the normalized statistical parameters such as rate of false positive, rate

308 of false negative, sensitivity, specificity, and accuracy.³¹ This suggested that the S/N ratio
 309 has a less relevant effect on the performance of this method compared to the F-ratio value.
 310 However, these results may be case dependent, therefore optimization of this parameter
 311 based on the data set should be considered by the analyst.

Table 2: The number of false positive detections, number of false negative detections, rate of false positive, rate of false negative, sensitivity, specificity, and accuracy parameters calculated for four different S/N values, having an F-ratio of 28, based on 20 simulations for each S/N.

Parameter	S/N values		
	1	3	10
False positive detection ^a (FP)	654	629	583
False negative detection ^a (FN)	122	151	166
Rate of false positive detection ^b (%)	9.0	9.0	8.0
Rate of false negative detection ^c (%)	2.0	2.5	2.7
Sensitivity ^e (%)	98.0	97.5	97.3
Specificity ^f (%)	90.7	91.1	91.7
Accuracy ^g (%)	94.1	94.1	94.3

^aThis parameter represents the number false positive detection out of total number of detections of 12,400, including 6,000 true positives (TP) and 6,400 true negatives (TN); ^bThe rate of false positive³¹ was calculated as FP/(FP+TN); ^cThe rate of false negative³¹ was calculated as: FN/(TP+FN); ^eThe sensitivity³¹ values were calculated using: TP/(TP+FN); ^fThe specificity³¹ values were calculated with: TN/(TN+FP); ^gThe accuracy³¹ values were calculated employing: (TP+TN)/(TP+FP+FN+TN).

312 Comparison between the unique feature list and the conventional 313 peak list

314 Once the F-ratio method was validated via semi-synthetic data, we processed the chro-
 315 matograms of the 15 sludge samples plus 3 method blanks using this algorithm. The same
 316 data set was also processed via MZmine, employing previously optimized parameters. The
 317 F-ratio method produced a list of unique features for each sample whereas MZmine created a
 318 conventional peak list for the same samples. We compared the unique feature lists produced
 319 via F-ratio method to the conventional peak lists by MZmine as well as the unique feature

320 lists produced via application of F-ratio method to both the whole chromatogram and the
321 peak list by MZmine. These comparisons enabled us to further evaluate/validate the F-ratio
322 method for analysis of the data generated via LC-HR-MS.

323

324 More than 92% of the unique features via F-ratio method were also present in the con-
325 ventional peak list via MZmine. For example, for one of the Oslo samples after the adducts
326 and isotopes removal 109 out of total 112 (i.e. 97%) unique features were also present in
327 the peak list of the same sample generated by MZmine. The number of features, via F-
328 ratio method, before adducts and isotope removal ranged from 403 features for one of the
329 blank samples to 127 for the Oslo sample whereas after the adducts and isotope removal
330 the unique features numbers ranged between 302 for the blank sample and 112 for the Oslo
331 sample. For the conventional peak list, we observed around 7500 peaks for the blank whereas
332 this number was around 12500 for the sludge samples. When comparing the unique feature
333 list to the conventional peak list, the number of discrepancy cases varied between 3 cases
334 for Oslo sample and 23 cases for the blank samples. A discrepancy case is defined as a
335 unique feature detected via F-ratio that is not present in the conventional peak list. All
336 the discrepancy cases were classified in two categories, for ease of explanation. The first
337 category and the most dominant one, particularly in the blank samples belonged to unique
338 features, which appeared to be noise rather than analytical signal. Considering the large
339 number of variables evaluated occurrence of a certain number of false positives was likely.
340 The second category was mainly caused by the fact that MZmine performs peak modeling
341 during the peak picking and uses the modeled apex for estimation of both m/z value and
342 the retention time. Using this approach, this algorithm may group shoulders of a peak with
343 the main peak. The F-ratio method however treats the shoulders as independent variables
344 and thus potential unique features. This category of discrepancy cases may also be caused
345 by the resolution of our instrument of 35,000. All considered, the F-ratio combined with the
346 apex detection algorithm method showed to have a large number (i.e. $\leq 92\%$) of common

347 unique features with the conventional peak picking approach which is an indication of its
348 robustness. Furthermore, these results imply that this method can be implemented in the
349 non-target workflows for structural elucidation.

350

351 The F-ratio applied to the whole chromatograms of the environmental samples resulted
352 in 250 unique feature in average while producing 3 unique features in average when applied
353 to the peak list generated via MZmine. A large number of the observed discrepancy cases
354 were due to the signal deconvolution, which was caused by the complexity of the analyzed
355 samples. The unsuccessful signal deconvolution was directly translated into the large within
356 group variability in the area of the integrated peaks, thus their lack of detection. The second
357 group of discrepancies was due to the peak modeling algorithm in MZmine, which failed to
358 detect the shoulder of a peak in the m/z domain as a separate peak, therefore their absence
359 from the unique feature list. It should be noted that the mentioned sources of failure in
360 the F-ratio applied to the peak list may be case dependent and may vary from dataset to
361 dataset. Further investigation of the potential sources of discrepancy between the F-ratio
362 variable selection applied to the whole chromatogram vs the peak list are needed.

363

364 The F-ratio method appeared to be able to successfully separate the sample groups
365 while both peak list and pixel-by-pixel methods failed in carrying out this task, Figure 1.
366 Multivariate statistical methods such as principal component analysis (PCA) when dealing
367 with large, complex datasets with a large level of noise and redundancy may fail to classify
368 the samples in logical groups. Consequently, univariate methods such as F-ratio are used
369 prior to these tests in order to reduce the redundancy in the data set. Therefore, a clear
370 and logical separation of the samples in the score plots is a crucial indication of a successful
371 prioritization/variable selection. We performed PCA on zero mask applied chromatograms
372 following the variable selection, the retention aligned peak list via MZmine, and the whole
373 chromatogram (i.e. pixel-by-pixel analysis). In the case of the sludge samples the inherent

374 complexity of the background signal was translated into inability of both peak list based
 375 and pixel-by-pixel based methods to separate these sample groups from each other properly.
 376 The F-ratio method, on the other hand, was able to perform separation of the sample groups
 377 because this method retains the variables that are causing the clustering of samples within a
 378 particular group. We also performed the F-ratio variable selection on the peak list generated
 379 via MZmine. In this case also the PCA was not able to separate the sample groups from
 380 each other, Figure 1. Therefore, it was not possible to perform a prioritization based on
 381 the peak list using the F-ratio method. Despite the mentioned complexity, the F-ratio
 382 method was able to separate the sample groups from each other, thus performing successful
 383 prioritization. These results also indicate the applicability of F-ratio method within the
 384 structure elucidation workflows during non-target analysis of complex samples analyzed via
 385 LC-HR-MS.

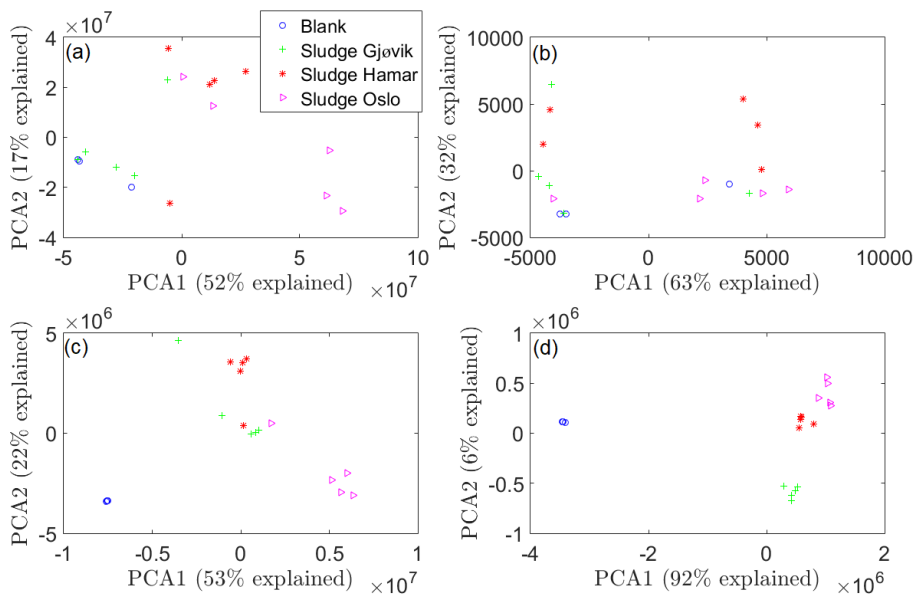


Figure 1: Figure depicting the PCA score plots of (a) peak list based classification, (b) peak list based after F-ratio variable selection (c) pixel-by-pixel analysis, and (d) the F-ratio method during the prioritization for non-target analysis of the 15 sludge samples plus 3 blanks.

Potential and Limitations

The F-ratio method combined with the apex detection showed to be a robust and reliable approach for prioritization of the unique features that are relevant to the sample classification. This method was effective at prioritization even for cases where the other conventional methods may fail due to the complexity of the analyzed data set, Figures 2 and 1. This method minimizes the data manipulations such as peak picking and/or modeling and at the same time results in a list of unique features which can be used for structure elucidation. The F-ratio method reduces the redundancy in the data set and detects the relevant variables in the data set enabling the analyst to focus on the identification of only the unique and relevant features. This method also has the advantage of being less dependent on the absolute intensity of the each chemical signal in the sample compared to the conventional prioritization methods. In other words, as long as a chemical signal causes large enough variability between the sample groups, independently from its absolute intensity, it will be detected as a unique relevant feature (Figure 2). Additionally, this method can be used for a battery of discovery-based non-targeted applications as long as there are replicates present. Furthermore, by changing the initial hypothesis, one can interrogate the data set in a completely different way. For example in case of the sludge samples in this study, by assuming that all the sludge samples belonged to one group and the blanks to another group, we could have selected the unique features which are in common in all the sludge samples and simultaneously subtracted the blanks from our samples.

There are also some limitations to application of F-ratio method for non-target analysis of LC-HR-MS data. This method is computationally expensive due to the large data sets produced when employing LC-HR-MS. For example, each sludge sample chromatogram in this study produced around 180 million variables (Figure 2), which requires a large computational power in order to be done in a timely manner. Moreover, this method has to be complemented with target and suspect analysis using the conventional methods for ubiq-

413 uitous chemicals or pollutants where measurable concentrations are more uniform. These
 414 pollutants would not be detected as unique features with statistically significant differences
 415 between sample groups. Also the use of data pre-treatment should be evaluated by the an-
 416 alyst on a case study base.

417

418 Considering capabilities and the limitations of the F-ratio method, this approach has a
 419 great potential to be applied to the LC-HR-MS non-target discovery-based analysis. The ap-
 420 plication of this method as well as its combination with the structural elucidation workflows
 421 are going to be subject of our future studies.

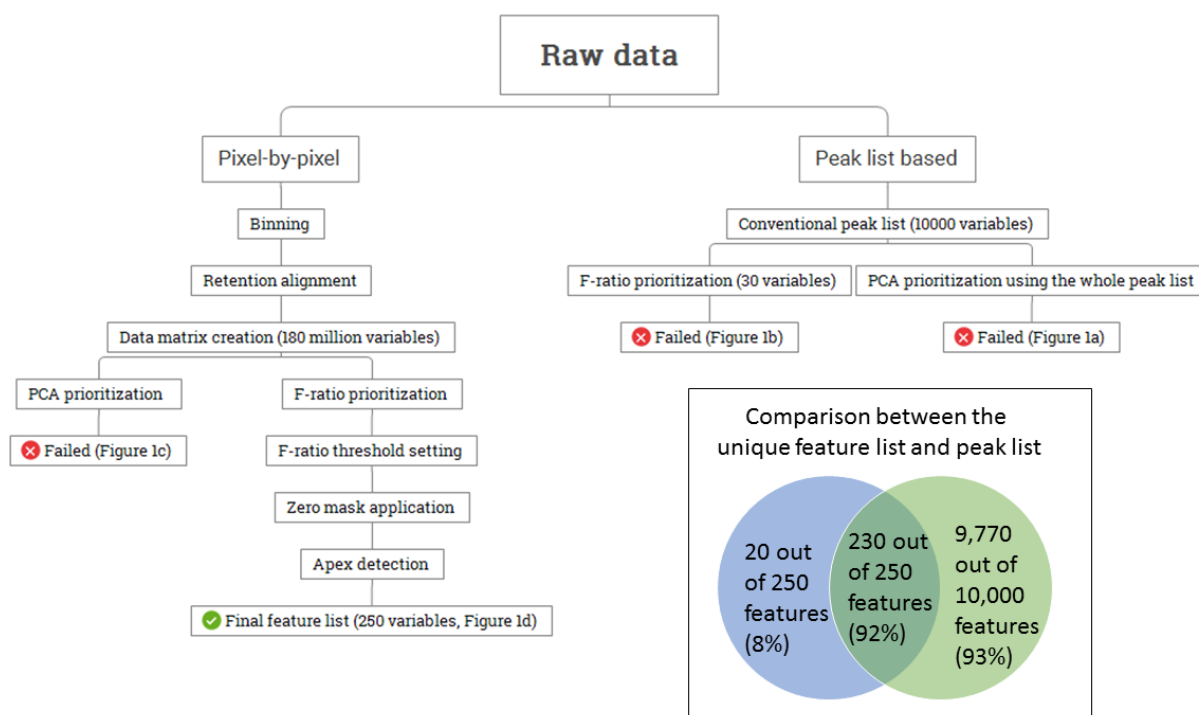


Figure 2: Figure depicting an overview of the F-ratio method vs the conventional methods as well as the venn diagrams of the comparison between the unique feature list and the conventional peak list generated via MZmine.

422 **Associated Content**

423 **Acknowledgement**

424 The authors are thankful to Prof. Bert van Bavel, Dr. Merete Grung and Dr. Jose A.
425 Baz-Lomba for their editorial input. We are also grateful to the Research Council of Norway
426 for the financial support of this project (RESOLVE, 243720).

427 **Supporting Information**

428 The Supporting Information including details regarding the sample preparation, analysis,
429 sateps taken during the data processing, and semi-synthetic data generation is available free
430 of charge on the ACS Publications website.

431 **Author Information**

432 Corresponding Author:

433 Saer Samanipour

434 E-mail: saer.samanipour@niva.no

435 Phone: +47 98 222 087

436 Address: Norwegian Institute for Water Research (NIVA)

437 0349 Oslo, Norway

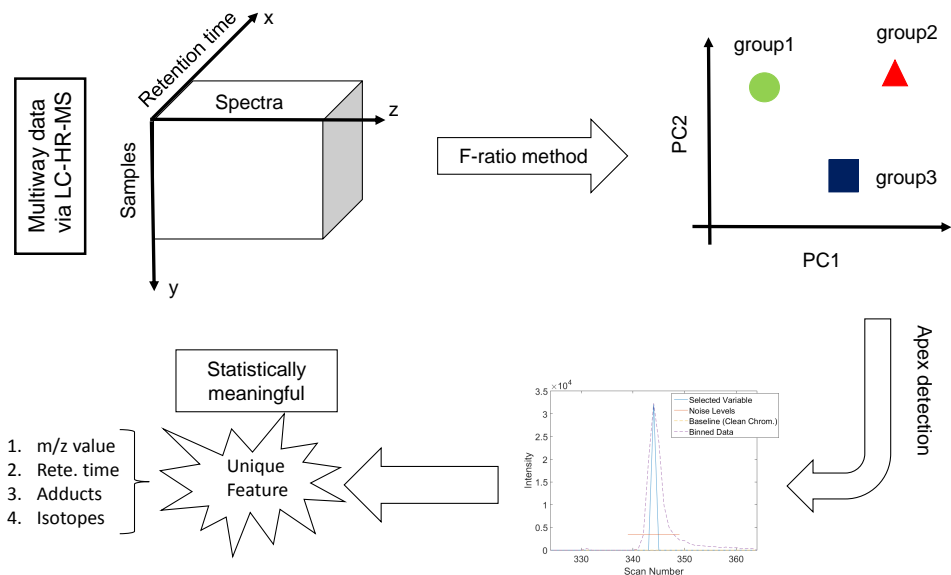
438 **References**

439 (1) Klečka, G.; Persoon, C.; Currie, R. *Reviews of Environmental Contamination and Tox-*
440 *icology Volume 207*; Springer, 2010; pp 1–93.

441 (2) La Farre, M.; Pérez, S.; Kantiani, L.; Barceló, D. *TrAC Trends Anal. Chem.* **2008**, *27*,
442 991–1007.

- 443 (3) Giger, W. *Anal. Bioanal. Chem.* **2009**, *393*, 37.
- 444 (4) Krauss, M.; Singer, H.; Hollender, J. *Anal. Bioanal. Chem.* **2010**, *397*, 943–951.
- 445 (5) Richardson, S. D. *Anal. Chem.* **2009**, *81*, 4645–4677.
- 446 (6) Schymanski, E. L.; Singer, H. P.; Longrée, P.; Loos, M.; Ruff, M.; Stravs, M. A.;
447 Ripollés Vidal, C.; Hollender, J. *Environ. Sci. Technol.* **2014**, *48*, 1811–1818.
- 448 (7) Gago-Ferrero, P.; Schymanski, E. L.; Bletsou, A. A.; Aalizadeh, R.; Hollender, J.;
449 Thomaidis, N. S. *Environ. Sci. Technol.* **2015**, *49*, 12333–12341.
- 450 (8) Samanipour, S.; Langford, K.; Reid, M. J.; Thomas, K. V. *J. Chromatogra. A* **2016**,
451 *1463*, 153–161.
- 452 (9) others,, et al. *Anal. Bioanal. Chem.* **2015**, *407*, 6237–6255.
- 453 (10) Gorrochategui, E.; Jaumot, J.; Lacorte, S.; Tauler, R. *Trends Anal. Chem.* **2016**, *82*,
454 425–442.
- 455 (11) Yi, L.; Dong, N.; Yun, Y.; Deng, B.; Ren, D.; Liu, S.; Liang, Y. *Anal. Chem. acta*
456 **2016**, *914*, 17–34.
- 457 (12) Chiaia-Hernandez, A. C.; Schymanski, E. L.; Kumar, P.; Singer, H. P.; Hollender, J.
458 *Anal. Bioanal. Chem.* **2014**, *406*, 7323–7335.
- 459 (13) Schollee, J. E.; Schymanski, E. L.; Avak, S. E.; Loos, M.; Hollender, J. *Anal. Chem.*
460 **2015**, *87*, 12121–12129.
- 461 (14) Thomas, K. V.; Langford, K.; Petersen, K.; Smith, A. J.; Tollefsen, K. E. *Environ. Sci.*
462 *Technol.* **2009**, *43*, 8066–8071.
- 463 (15) Kalogiouri, N. P.; Alygizakis, N. A.; Aalizadeh, R.; Thomaidis, N. S. *Anal. and Bioanal.*
464 *Chem.* **2016**, *408*, 7955–7970.

- 465 (16) Pierce, K. M.; Hoggard, J. C.; Hope, J. L.; Rainey, P. M.; Hoofnagle, A. N.; Jack, R. M.;
466 Wright, B. W.; Synovec, R. E. *Anal. Chem.* **2006**, *78*, 5068–5075.
- 467 (17) Beckstrom, A. C.; Humston, E. M.; Snyder, L. R.; Synovec, R. E.; Juul, S. E. *J.*
468 *Chromatogr. A* **2011**, *1218*, 1899–1906.
- 469 (18) Christensen, J. H.; Tomasi, G. *J. Chromatogr. A* **2007**, *1169*, 1–22.
- 470 (19) Parsons, B. A.; Marney, L. C.; Siegler, W. C.; Hoggard, J. C.; Wright, B. W.; Syn-
471 ovec, R. E. *Analytical chemistry* **2015**, *87*, 3812–3819.
- 472 (20) Parsons, B. A.; Pinkerton, D. K.; Wright, B. W.; Synovec, R. E. *J. Chromatogr. A*
473 **2016**, *1440*, 179–190.
- 474 (21) Sinkov, N. A.; Sandercock, P. M. L.; Harynuk, J. J. *Forensic Sci. Int.* **2014**, *235*, 24–31.
- 475 (22) Sinkov, N. A.; Harynuk, J. J. *Talanta* **2011**, *83*, 1079–1087.
- 476 (23) Sinkov, N. A.; Harynuk, J. J. *Talanta* **2013**, *103*, 252–259.
- 477 (24) Adutwum, L.; Harynuk, J. *Anal. Chem.* **2014**, *86*, 7726–7733.
- 478 (25) Watson, N. E.; VanWingerden, M. M.; Pierce, K. M.; Wright, B. W.; Synovec, R. E.
479 *J. Chromatogr. A* **2006**, *1129*, 111–118.
- 480 (26) Baz-Lomba, J. A.; Reid, M. J.; Thomas, K. V. *Anal. Chem. acta* **2016**, *914*, 81–90.
- 481 (27) MATLAB version 9.1 Natick, Massachusetts: The MathWorks Inc.,
- 482 (28) Nadeau, J. S.; Wright, B. W.; Synovec, R. E. *Talanta* **2010**, *81*, 120–128.
- 483 (29) Katajamaa, M.; Miettinen, J.; Orešič, M. *Bioinformatics* **2006**, *22*, 634–636.
- 484 (30) Brereton, R. G. *Applied chemometrics for scientists*; John Wiley & Sons, 2007.
- 485 (31) Burke, D. S.; Brundage, J. F.; Redfield, R. R.; Damato, J. J.; Schable, C. A.; Put-
486 man, P.; Visintine, R.; Kim, H. I. *N. Engl. J. Med.* **1988**, *319*, 961–964.



for TOC Only