

1 **TITLE:** Rapid and accurate identification of *Escherichia coli* STEC O157:H7 by  
2 mass spectrometry, artificial intelligence and detection of specific biomarkers  
3 peaks.

4  
5 **Authors:** Manfredi Eduardo <sup>a,1,\*</sup>, Rocca María Florencia <sup>b,c,1,\*</sup>, Zintgraff Jonathan <sup>b,c,1,\*</sup>, Irazu  
6 Lucía <sup>b</sup>, Miliwebsky Elizabeth <sup>a</sup>, Carbonari Carolina <sup>a</sup>, Deza Natalia <sup>a</sup>, Prieto Monica <sup>b,c</sup>, Chinen  
7 Isabel <sup>a</sup>

8  
9  
10  
11  
12  
13

14 <sup>a-</sup> *Servicio Fisiopatogenia, Instituto Nacional de Enfermedades Infecciosas (INEI) –*  
15 *Administración Nacional de Laboratorios e Institutos de Salud (ANLIS) “Dr. Carlos G.*  
16 *Malbrán”*

17  
18  
19  
20

18 <sup>b-</sup> *Instituto Nacional de Enfermedades Infecciosas (INEI) – Administración Nacional de*  
19 *Laboratorios e Institutos de Salud (ANLIS) “Dr. Carlos G. Malbrán”*

21  
22  
23

21 <sup>c-</sup> *Red Nacional de Espectrometría de Masas aplicada a la Microbiología Clínica (ReNaEM*  
22 *Argentina), Argentina*

24  
25  
26  
27  
28  
29

30 **\*Corresponding author at:** *Instituto Nacional de Enfermedades Infecciosas (INEI) –*  
31 *Administración Nacional de Laboratorios e Institutos de Salud (ANLIS) “Dr. Carlos G. Malbrán*  
32

33  
34  
35  
36  
37

**E-mail addresses:**

34 [edimanfredi@gmail.com](mailto:edimanfredi@gmail.com) (E.M)  
35 [florirocca1980@gmail.com](mailto:florirocca1980@gmail.com) (M.F.R)  
36 [jczintgraff@gmail.com](mailto:jczintgraff@gmail.com) (J.Z)

38  
39

1<sup>-</sup> These authors contributed equally to this work.

40  
41  
42  
43  
44  
45

46  
47  
48

## ABSTRACT

49 The different pathotypes of *Escherichia* can produce a large number of human  
50 diseases. Surveillance becomes complex since their differentiation are not  
51 easy.

52 Particularly, the detection of Shiga toxin-producing *Escherichia coli* (STEC)  
53 serotype O157:H7 consists of stool culture of a diarrheal sample in enrichment  
54 and/or selective media, identification of presumptive colonies and confirmation  
55 by Multiplex PCR technique for the genotypic characterization of serogroup  
56 O157 and Shiga toxins (*stx1* and *stx2*), in addition to the traditional biochemical  
57 identification.

58 All of these procedures are laborious, require a certain level of training, are time  
59 consuming and expensive. Among the currently most widely used  
60 methodologies, MALDI-TOF MS mass spectrometry (matrix-assisted laser  
61 desorption/ionization with time-of-flight mass detection), allows a quick and  
62 easy way to obtain a protein spectrum of a microorganism, not only in order to  
63 identify the genus and species, but also the discovery of potential biomarker  
64 peaks of a certain characteristic. In the present work, the information obtained  
65 from 60 clinical isolates was used to detect peptide fingerprints of STEC  
66 O157:H7 and other diarrheagenic *E. coli*. The differences found in the protein  
67 profiles of the different pathotypes established the foundations for the  
68 development and evaluation of classification models through automated  
69 training.

70 The application of the Biomarkers in combination with the predictive models on  
71 a new set of samples (n=142), achieved 99.3% of correct classifications,

72 allowing the distinction between STEC O157:H7 isolates from the other  
73 diarrheal *Escherichia coli*.

74 Therefore, given that STEC O157:H7 is the main causal agent of haemolytic  
75 uremic syndrome and based on the performance values obtained in the present  
76 work (Sensitivity=98.5% and Specificity=100%), this development could be a  
77 useful tool for diagnosis of the disease in clinical microbiology laboratories.

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99 **INTRODUCTION.**

100

101 The contribution of *Escherichia coli* to human intestinal disease may be largely  
102 uncharacterized, because many types of pathogenic *E. coli* are not routinely  
103 tested in clinical microbiology laboratories.

104 Shiga toxin-producing *Escherichia coli* (STEC) is associated with outbreaks that  
105 causes diarrhea, hemorrhagic colitis, and hemolytic-uremic syndrome (HUS) in  
106 humans. It is part of the diarrheagenic *E. coli* (DEC) group, which also includes:  
107 enteropathogenic *E. coli* (EPEC), enterotoxigenic *E. coli* (ETEC),  
108 enteroaggregative *E. coli* (EAEC), enteroinvasive *E. coli* (EIEC) and diffusely  
109 adherent *E. coli* (DAEC). Although, there are more than 150 serotypes [1] that  
110 share the same pathogenic potential, O157:H7 is the most frequent. In  
111 particular, the detection of STEC O157:H7 consists of the culture of the faecal  
112 sample in enrichment and/or selective media such as MacConkey agar with  
113 sorbitol for the identification of presumptive non-fermenting sorbitol colonies and  
114 confirmation by Multiplex PCR for the genotypic characterization of serogroup  
115 O157 and Shiga toxins (*stx1* and *stx2*), in addition to subsequent traditional  
116 biochemical identification.

117 However, all this methodological complexity is very difficult to implement in a  
118 traditional laboratory [2].

119 On the other hand, mass spectrometry (MS), specifically MALDI-TOF MS  
120 (matrix-assisted laser desorption/ionization with time-of-flight mass detection),  
121 provides a simple, rapid, robust, and low-cost microbial identification. MALDI-

122 TOF MS is a technique based on the analysis of protein spectra containing  
123 peaks with an exactly determinable mass-charge ratio ( $m/z$ ) generated by the  
124 impact of a laser on a previously crystallized isolate with an organic matrix. In  
125 recent years, MS has acquired great importance in the identification of  
126 pathogens that are clinically relevant in public health [3-5]. However, the  
127 potential of this methodology combined with machine learning algorithms for the  
128 detection of profiles in a wide variety of samples and its use as a screening  
129 technique is expanding, due to its low-cost and high performance [6].

130 In this study, we wanted to verify the usefulness of MS to rapidly identify  
131 O157:H7 STEC from pathotypes other than diarrheagenic *E. coli*; then, we  
132 proposed to detect and analyse peaks in the spectra generated by MALDI-TOF  
133 MS to find possible biomarkers and thus establish differential patterns between  
134 a wide variety of *E. coli* strains.

135

## 136 **MATERIALS AND METHODS.**

137

### 138 **Isolates collection.**

139 The spectra obtained from 60 isolates corresponding to four different DEC  
140 categories were used for the development of predictive models and the  
141 detection of biomarkers: EPEC (n=15), ETEC (n=15), STEC NON O157 (n=20)  
142 and STEC O157:H7 (n=10). The detail of the isolates can be found in **Table S1**  
143 in Supplementary Material.

144 For the final validation, we used 142 different isolates of: STEC O157:H7  
145 (n=65), non-toxicogenic *E.coli* O157 (n=13), STEC NON O157 (n=17), ETEC (n=  
146 11), EPEC (n=12), EAEC (n=15), EIEC (n=7), and *E. coli* without virulence

147 factors (n=2). All of them were obtained mainly from faecal samples from  
148 different health institutions in our country and food samples subsequently  
149 referred to the National Reference Laboratory- Servicio de Fisiopatogenia INEI-  
150 ANLIS "Dr. Carlos G. Malbrán"- for confirmation and characterization of specific  
151 virulence factors [7].

152

### 153 **Acquisition of spectra.**

154 The Microflex LT mass spectrometry equipment (Bruker Daltonics) was used to  
155 obtain the protein spectra from *E. coli* isolates. Subsequently, each isolate was  
156 spotted in quadruplicate in the wells of the steel plate provided by the  
157 manufacturer using the direct method and crystallized by adding 1 ul of HCCA  
158 matrix ( $\alpha$ -Cyano-4-hydroxycinnamic acid in 50% of acetonitrile and 2.5%  
159 trifluoroacetic acid). After a few minutes of drying, the plate was introduced into  
160 the equipment and once the vacuum was reached in the Flex Control v3.4  
161 software, the spectra were acquired in the linear positive mode, with 30-40%  
162 laser power. and in a mass range of 2 to 20 kDa. Each well was read twice, so  
163 eight individual spectra were obtained for each isolate, thus minimizing the  
164 variability of the technique.

165 The external calibrator provided by the manufacturer, BTS (Bruker Test  
166 Standard), was used prior to each run.

167 The 142 isolates used in the subsequent validation set were processed in the  
168 same way by direct method, but each isolate was spotted in duplicate and read  
169 only once, simulating the routine procedure of a microbiology laboratory.

170

171

172

173

174

175 **MALDI-TOF analysis**

176 All isolates were identified using the MALDI Biotyper RTC software. According  
177 to the manufacturer's recommendations, the identification is considered reliable  
178 at the species level when the score values greater than 2.0 are obtained. When  
179 the score values are between 1.7 and 1.99, it is considered reliable  
180 identification at the genus level; and it is considered 'No Identification' when the  
181 value of the score is  $\leq 1.69$  [8].

182

183 **Bioinformatic analysis.**

184 To perform data analysis, ClinPro Tools software (version 3.0, Bruker Daltonik  
185 GmbH, Bremen, Germany) and Flex Analysis v3.4 software (Bruker Daltonics,  
186 Bremen, Germany) were used.

187

188 **Data pre-processing.**

189 In order to take advantage of the greatest amount of information contained in  
190 the spectra, the following data pre-processing steps were performed: baseline  
191 correction (top hat 10% of minimum width of the baseline), smoothing and  
192 calibration excluding null spectra or out of range, according to the literature [9-  
193 11].

194

195 **Peak Selection.**

196 The exploration of the potential biomarkers was performed on the protein  
197 profiles generated from the 60 isolates that were part of the initial training group  
198 and which were also, used to create the classification models.

199

#### 200 **Flex Analysis v3.4 Software.**

201 All spectra were exported as mzXML files using CompasXport CXP3.0.5 and a  
202 series of analyses were performed according to standard Bruker setup.  
203 Spectrum quality criteria for overall aspect and intensity were checked. Next,  
204 visually identifiable biomarker peaks were explored in the different views offered  
205 by the program. The mass list was exported to Excel (Microsoft, Redmond, WA)  
206 to analyse possible biomarker peaks. Values of "1" or "0" (data binarization)  
207 were assigned to the presence or absence of a peak within the tolerance  
208 interval (+/- 10Da). Based on this analysis, groups of potentially useful peaks for  
209 the diagnosis of STEC O157:H7 were found.

210

#### 211 **ClinProTools v3.0 Software**

212 The spectra of STEC O157:H7 were assigned as class 1 and the rest of the  
213 DEC isolates were class 2. Biomarker peaks were automatically identified by  
214 class comparison using the function "Peak Statistic Table".

215 To select the characteristic peaks of the two classes, the following statistical  
216 tests were used: the t test/analysis of variance ANOVA (PTTA), Wilcoxon or  
217 Kruskal–Wallis test (W/KW), and Anderson–Darling test (AD). A p value of 0.05  
218 was established as the cut-off point **[12]**:

219 -if p is <0.05 in the AD test, a characteristic peak is selected if the  
220 corresponding p-value in the W/KW test is also <0.05.



221 -if p is 0.05 in the AD test, then a characteristic peak is selected if the  
222 corresponding p value in ANOVA is also <0.05 [13].

223 The discriminative power for each biomarker was further described by receiver  
224 operating characteristic (ROC) and area under the curve (AUC) analysis.

225 ROC curve indicates the relationship of the true-positive rate (TPR) and the  
226 false-positive rate (FPR). The area under the ROC curve is equal to the  
227 probability that a biomarker sorts a randomly selected positive sample higher  
228 than a randomly selected negative one.

229

### 230 **Principal Component Analysis (PCA).**

231 To explore and compare spectra in multidimensional space and in order to  
232 evaluate the possible distributions or clusters on the isolates of both classes, a  
233 first exploratory and unsupervised analysis was performed of all 60 samples.

234

### 235 **Classification models.**

236 Supervised classification models were performed using the following algorithms  
237 provided by ClinPro Tools software: Supervised Neural Network (SNN),  
238 optimized genetic algorithm combined with k-nearest neighbour classification  
239 (GA/ kNN) and a quickclassifier (QC).

240 For each model, the following parameters were calculated: Recognition  
241 Capability (CR) and Cross Validation Percentage (VC), both of which are  
242 indicators of the theoretical behaviour that the model will have in future  
243 classifications.

244

### 245 **Selection of isolates for final validation.**

246 To evaluate the robustness of the models created, an independent set of  
247 isolates different from those used to developed the algorithms was selected  
248 (N=142). For each isolate, a spectrum was presented to the selected  
249 classification model. The software then returned a result that was compared to  
250 current reference techniques.

251

### 252 **Statistical analysis.**

253 The parameters evaluated were: sensitivity, specificity, positive predictive value,  
254 negative predictive value [14]. Besides, the CLSI guide, EP12-A2, was used to  
255 compare methods that report results qualitatively. When the comparison is  
256 made with a method that is not considered a reference, the degree of similarity  
257 between the methods is measured through the percentage of negative  
258 agreement and the percentage of positive agreement. The diagnostic  
259 parameters of the methods are then compared to determine if the difference  
260 between the two of them is statistically significant.

261

## 262 **RESULTS.**

263

### 264 **Confirmation at the genus-species level**

265 All isolates were identified at the species level as *Escherichia coli* with a score  
266 value greater than 2.0, in agreement with the result of the reference techniques.

267

### 268 **Unsupervised Analysis.**

269 First, a Principal Component Analysis (PCA) was performed for the 60 isolates  
270 used as a training set. In this way it was possible to reduce all the information  
271 contained in the MALDI-TOF spectra in a few new variables.

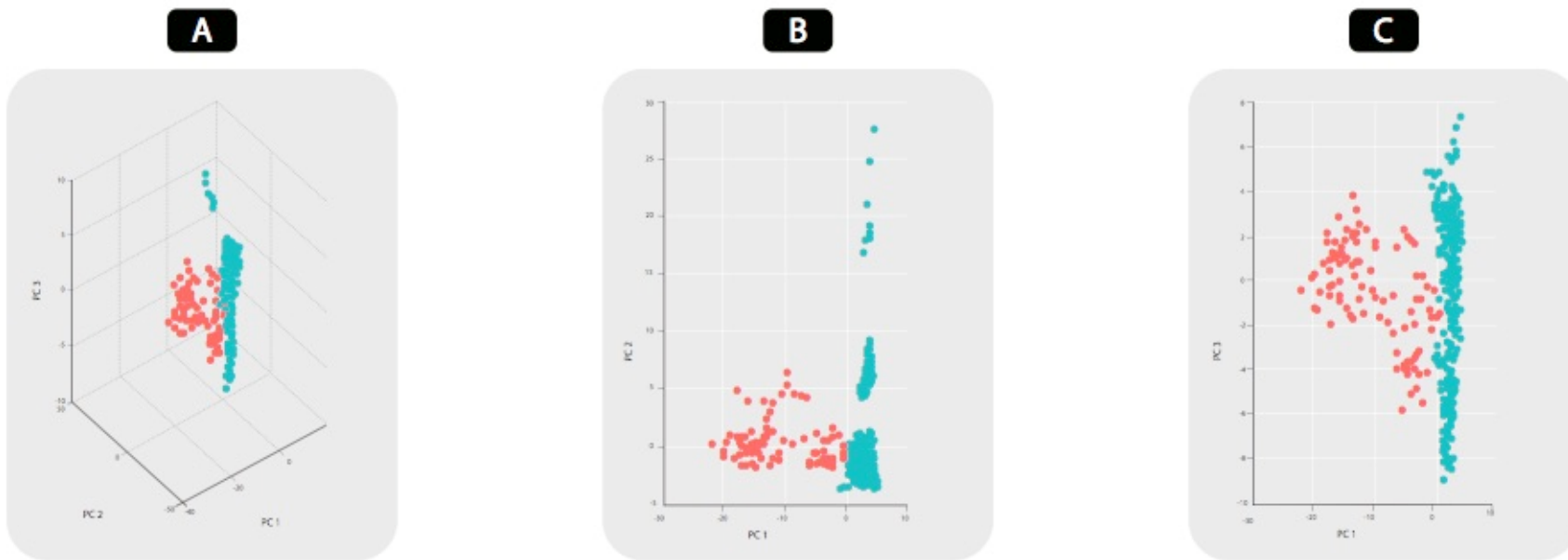
272 This allows us to graphically represent all the spectra together in three and two  
273 dimensions on the Score plot. (Figure 1) The complete list of significant peaks  
274 found in ClinPro Tools can be found in **Supplementary Material S2**.

275

276

277

278



·DEC  
·STEC O157:H

280 **Figure 1.** PCA plots results. Data originated from the external MATLAB software tool integrated into ClinPro Tools.

281 Graph A shows the three components plotted simultaneously in three  
282 dimensions, while in graphs B and C shows PC1 versus PC2 and PC1 versus  
283 PC3 respectively.

284

### 285 **Supervised Analysis.**

286 Subsequently, a supervised multivariate analysis was performed with the  
287 additional information of each isolate to define each class:

288

289 **CLASS 1: STEC O157:H7**

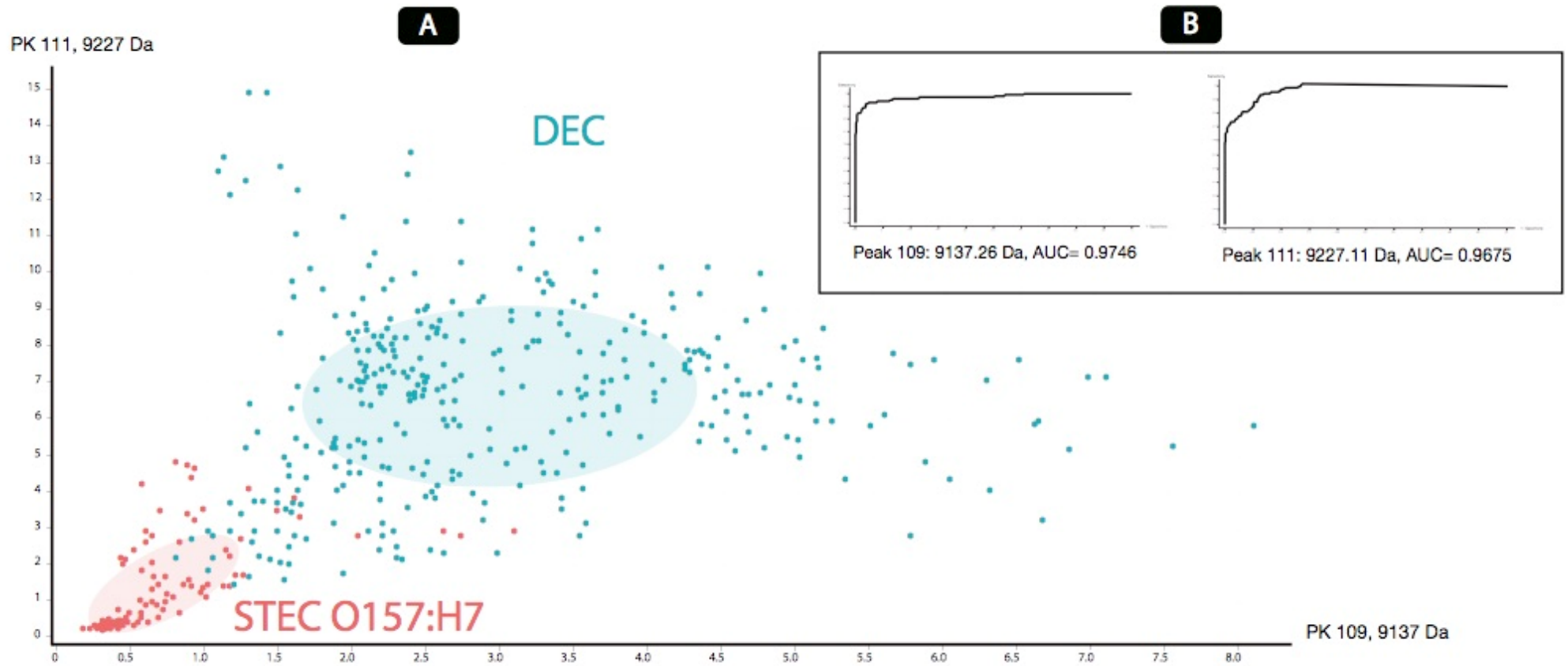
290 **CLASS 2: NON STEC O157 (other DEC)**

291

292 Figure 2-A shows the two-dimensional distribution plot of all the spectra of each  
293 class based on the two best peaks obtained for their classification; which  
294 correspond to the 9137.26 Da peak and the 9227.11 Da peak. The peak  
295 number and its m/z values are shown on the x and y axes respectively, while  
296 the ellipses represent the 95% confidence interval. On the other hand, Figure 2-  
297 B shows the ROC curves of the two selected peaks. The area under the curve  
298 (AUC) represents the discriminatory potential of each biomarker peak.

299

300



302 **Figure 2.** **A-** 2D graph of peak distribution of the 2-class model. **B-** ROC curves of the most discriminating peaks according to the analyses performed

303

304 **Classifier Models.**

305 The predictive models were calculated based on three available algorithms:

306 GA/kNN, SNN and QC, the results of the different parameters of each algorithm

307 are summarized in **Table 1**.

308 The SNN algorithm was discarded from the successive analyses due that the

309 results obtained were not optimal.

310

311 **Table 1.** Results of RC, CV and integration areas for each algorithm.

312

| Classifiers<br>Algorithms | RC      | CV      | Integration areas used by each<br>model |
|---------------------------|---------|---------|---|
| <b>GA / kNN</b>           | 100.00% | 100.00% | 3082;4939;5080;8813;8994                |
| <b>QC</b>                 | 92.83%  | 92.47%  | 6389;9136;9226                          |

RC=  
Rec  
ognit  
ion

317 Capability; VC=Cross Validation

318

319 As a result of the external validation carried out with the set of 142 isolates, a

320 good performance was observed with the GA/kNN algorithm and with the QC

321 algorithm. Nevertheless, we decided to combine both models; first, using the

322 QC algorithm but applying a cut-off value  $\geq 1.55$ , since from this value 100%

323 correct classifications were observed compared to the reference technique

324 (Figure 3).

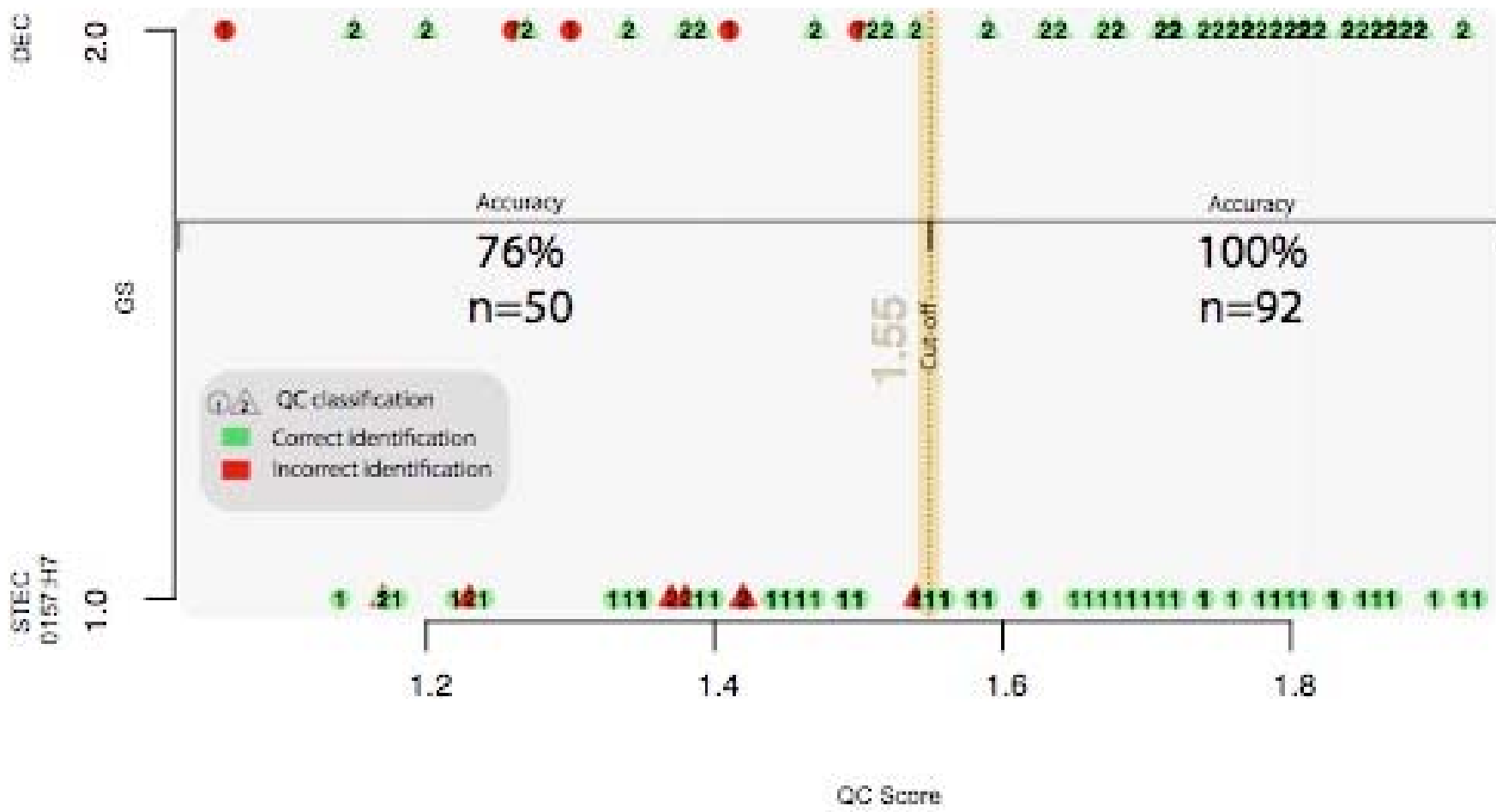
325

326

327

328

329



3:

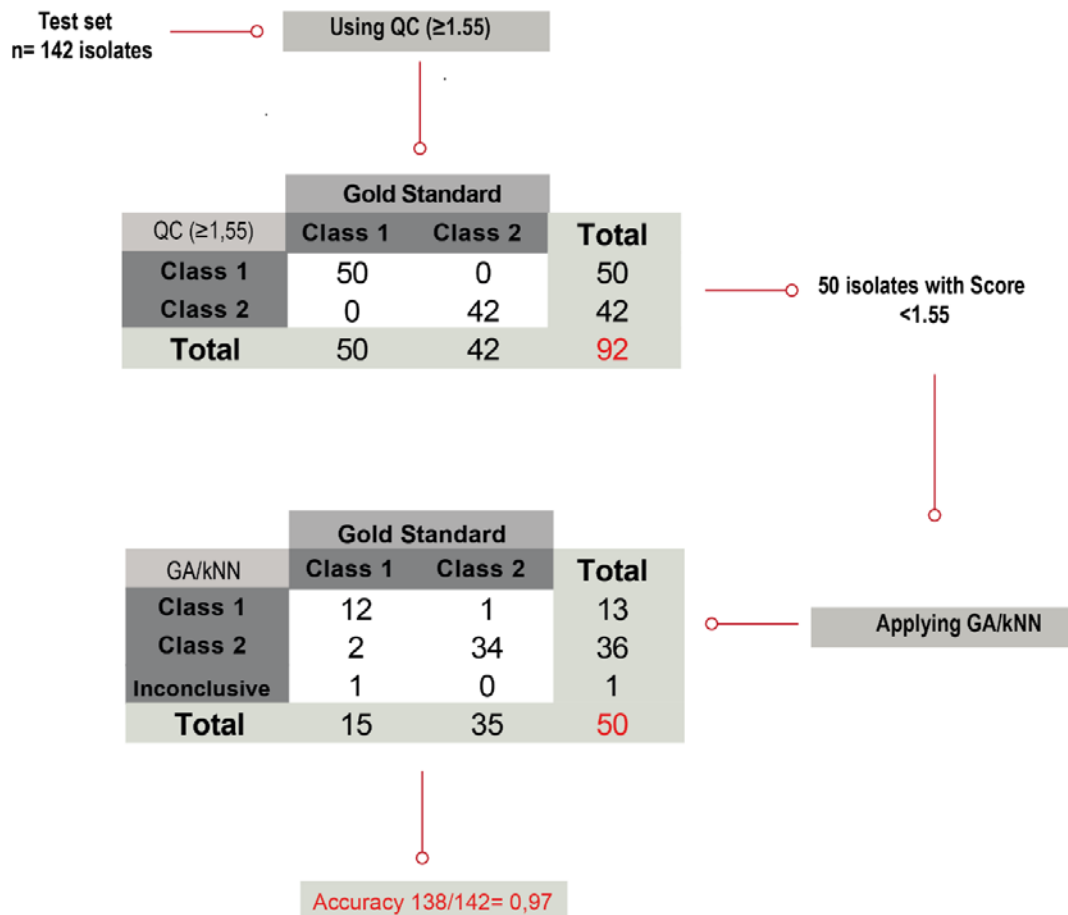
331 **Figure 3.** Graph of QC scores where the 100% concordance of the QC algorithm with the reference method is observed from the standardized cut-off value of  
 332 1.55.



333 On the other hand, the isolates that were classified with a QC value  $<1.55$  the  
334 GA/kNN algorithm was applied (Figure 4). This combination managed to  
335 increase the identification up to 97%, as detailed later in Table 3.

336 There were three isolates incorrectly classified using this scheme and one was  
337 considered Inconclusive since the result of the QC algorithm was  $<1.55$  and  
338 when applying the GA/kNN algorithm dissimilar results were obtained between  
339 the sample and the duplicate.

340



341

342 **Figure 4.** Algorithm applied for the identification of STEC O157:H7 based on predictive  
343 classification models.

344

345

346 **Biomarker detection.**

347 10 differential peaks with statistical significance were found in both software; of  
348 which 9 correspond to the STEC O157:H7 pathotype and a single biomarker  
349 was present only in NON STEC O157 (Figure 5).

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

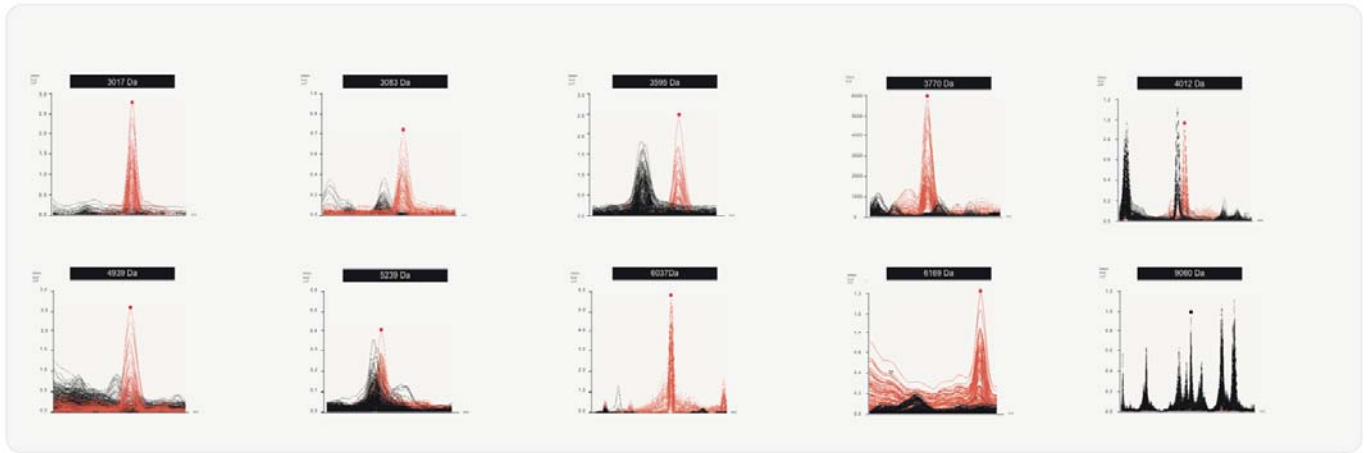
367

368

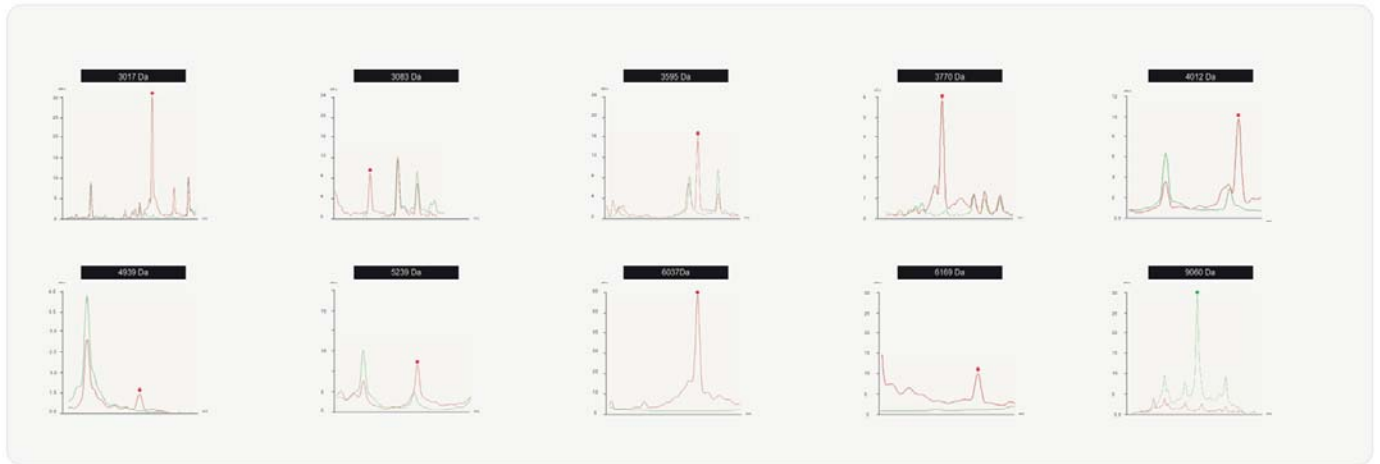
369

370

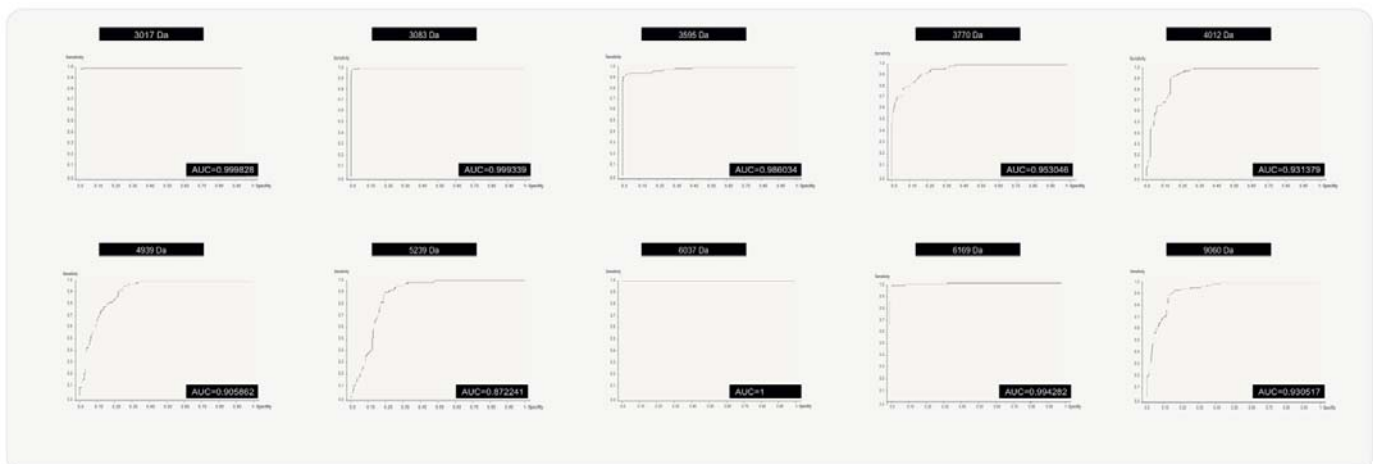
**A**



**B**



**C**



372 **Figure. 5. A-** Characteristic peaks (biomarkers) in individual spectra of STEC O157:H7 samples  
 373 (red) versus DEC samples (black), obtained by manual analysis in Flex Analysis v3.4 software.  
 374 **B-** Average spectra of the same peaks, STEC O157:H7 (red); DEC (green), obtained by ClinPro  
 375 Tools v3.0. **C-** ROC curves and AUC values originated from the external MATLAB software tool  
 376 integrated in ClinPro Tools.

377

378 The profile of the 10 potential biomarkers selected for the differentiation of  
 379 STEC O157:H7 from DEC isolates is shown in Table 2. The description of the  
 380 profiles found for all the challenged isolates is found in Table S3 of the  
 381 Supplementary Material.

382

383 **Table 2.** Profile of the 10 potential biomarkers selected for the differentiation of STEC O157:H7  
 384 from DEC.

| Classification                    | Biomarkers (m/z) |                |               |               |                |                |                |               |                |                |
|-----------------------------------|------------------|----------------|---------------|---------------|----------------|----------------|----------------|---------------|----------------|----------------|
|                                   | 3017             | 3083           | 3595          | 3770          | 4012           | 4939           | 5238           | 6037          | 6169           | 9060           |
| <b>CLASS 1</b><br>(O157:H7; n=65) | 72%<br>(n=47)    | 45%<br>(n=29)  | 51%<br>(n=33) | 85%<br>(n=55) | 40%<br>(n=26)  | 38%<br>(n=25)  | 97%<br>(n=63)  | 86%<br>(n=56) | 55%<br>(n=36)  | 100%<br>(n=65) |
| <b>CLASS 2</b><br>(DEC; n=77)     | 96%<br>(n=74)    | 100%<br>(n=77) | 99%<br>(n=76) | 99%<br>(n=76) | 100%<br>(n=77) | 100%<br>(n=77) | 100%<br>(n=77) | 97%<br>(n=75) | 100%<br>(n=77) | 96%<br>(n=74)  |

385

386 Based on the analysis of the results obtained from the BM challenge on 142  
 387 new samples in duplicate, it can be confirmed that the absence of the 9060 Da  
 388 peak, added to the detection of at least one of the other nine peaks described in  
 389 this work, confirmed the identification of an isolate of STEC O157:H7, although  
 390 most of these isolates (92%), in addition to not showing the peak at 9060 Da,  
 391 had 3 or more biomarkers.

392 One isolate of STEC O157:H7 and 3 other isolates (STEC NON O157, ETEC,  
 393 and non-toxigenic O157) did not present any of the peaks listed above, and  
 394 thus could not be classified.

395 Regarding the NON STEC O157, 96% presented the peak of m/z 9060 Da and  
 396 93.5% none of the nine peaks previously described.

397 Finally, the sensitivity of the search for probable biomarker peaks was 96.9%  
 398 and the specificity 100%, as shown in Table 3.

399

400 **Table 3.** Sensitivity, Specificity, PPV and NPV values of the different approaches evaluated and  
 401 the statistical relationship between them.

| Evaluated<br>Parameter | Approaches      |             |                                  |            |   |
|------------------------|-----------------|-------------|----------------------------------|------------|---|
|                        | A               | B           | C                                | D          | E   |
|                        | Model<br>GA/kNN | Model<br>QC | Combination<br>QC<br>+<br>GA/kNN | Biomarkers | Combination<br>QC + GA/kNN<br>+<br>Biomarkers |
| <b>Sensitivity</b>     | 80,00%          | 92,30%      | 95,40%                           | 96,90%     | 98,50%  |
| <b>Specificity</b>     | 96,10%          | 90,90%      | 98,70%                           | 100%       | 100%  |
| <b>PPV</b>             | 94,50%          | 89,60%      | 98,40%                           | 100%       | 100%  |
| <b>NPV</b>             | 85,10%          | 93,30%      | 96,20%                           | 97,50%     | 98,70%  |

402 **PPV=** Positive predictive value; **NPV=** Negative predictive value

403

404 The difference between the sensitivity and specificity of the mathematical model  
 405 with respect to the manual (C and D) and the two best methods (D and E) and  
 406 with respect to the reference method was estimated to conclude whether these  
 407 differences were statistically significant or not. In this way, if the 95% confidence  
 408 interval of the differences contains the value 0, it is concluded that there are no

409 statistically significant differences, otherwise there are differences. The results  
410 of this analysis are detailed in Table 4.

411

412 **Table 4.** Results of the difference between the sensitivity and specificity of the mathematical  
413 and manual model and of the two best methods with respect to the reference method.

| Combination<br>of approaches | Difference<br>sensitivity (%) | Difference<br>specificity (%) | Confidence<br>Intervals |       | Conclusion  |
|------------------------------|-------------------------------|-------------------------------|-------------------------|-------|---|
|                              |                               |                               | 95%                     |       |   |
|                              |                               |                               | Lower                   | Upper |   |
| C vs D                       | -1.54                         | -                             | -8.58                   | 6.1   | The difference between the sensitivity of the methods is <b>not statistically significant</b> |
|                              | -                             | 1.3                           | -3.57                   | 7.00  | The difference between the sensitivity of the methods is <b>not statistically significant</b> |
| D vs E                       | 1.56                          | -                             | -4.52                   | 9.03  | The difference between the sensitivity of the methods is <b>not statistically significant</b> |
|                              | -                             | 0                             | -5.63                   | 4.87  | The difference between the sensitivity of the methods is <b>not statistically significant</b> |

414

415

416 Therefore, despite the fact that there were no statistically significant differences  
417 between the performance values of approaches D and E, the benefits of their  
418 combined use are evident, as described in numerous publications in the current  
419 literature **[15]** and it follows from this work in the resolution of discordant cases.

420

421 In summary, if both developments are applied in a complementary way to  
422 isolates that could not be correctly classified by automated training or did not  
423 present any of the peaks considered biomarkers, accurate detection of 98.5% of  
424 STEC O157:H7 isolates is achieved and the correct classification of 99.3% of all  
425 the isolates studied

426

427 It was observed that 3/3 isolates incorrectly classified by the predictive models  
428 were correctly resolved by the BM finding method and 3/4 that could not be  
429 classified because they did not present the peaks, were resolved by  
430 mathematical models; a single case could not be resolved by either of the two  
431 methods, reaffirming the usefulness of the combined use of both approaches  
432 (Table 5).

433 The table of the results obtained on the total number of isolates applying  
434 machine learning and biomarker detection, in comparison with current reference  
435 techniques, can be found in Table S3 the Supplementary Material.

436

437 **Table 5.** Discordant cases of the different approaches compared with the gold standard results.

| Samples ID | GA/kNN+QC | Biomarkers | GA/kNN + QC<br>+<br>Biomarkers | Gold<br>standard |
|------------|-----------|------------|--------------------------------|------------------|
| 750/18     | Class 2   | Class 1    | Class 1                        | Class 1          |
| 506/18     | Class 2   | N/D        | Class 2                        | Class 2          |
| 385/16     | Class 1   | Class 2    | Class 2                        | Class 2          |
| 504/18     | Class 2   | Class 1    | Class 1                        | Class 1          |
| 329/18     | Class 2   | N/D        | Class 2                        | Class 2          |
| 714/18     | N/D       | N/D        | N/D                            | Class 1          |
| 493/18     | Class 2   | N/D        | Class 2                        | Class 2          |

438 N/D=not determinated

439

#### 440 **CONCLUSIONS.**

441 Due to the important analytical capabilities that MS currently has, added to the  
442 speed of results and lower-cost, the possible implementation of the MALDI-TOF  
443 MS system coupled to simple and practical artificial intelligence tools could be  
444 considered as a STEC O157:H7 diagnostic screening method.

445 Through the proteomic analysis of the information contained in the spectra of  
446 the different classes of *E. coli*, and applying a combination of predictive models  
447 based on machine learning, it was possible to quickly identify 94% of the STEC  
448 O157:H7 isolates and precise, starting from characteristic suspicious colonies,  
449 which implied a substantial saving of time and resources in the routine of the  
450 conventional laboratory. By combining this approach with the search for  
451 potential biomarker peaks, the percentage of correct identifications rose to  
452 98.5%.

453 There were several previous attempts in the literature to detect STEC O157 by  
454 MALDI-TOF Mass Spectrometry [16-19], however, no defining peaks were  
455 found in any of the previous works. and without the requirement of complex  
456 extraction techniques or equipment with greater discriminatory power, such as  
457 the TOF-TOF type, or peak readings above 10,000 Da, which are generally less  
458 detected. On the other hand, here we detect a large number of reproducible  
459 peaks in the reading range of the order of 3000 to 9000Da by direct method,  
460 without the need to make any modifications, which results in a simple, fast and  
461 easily transferable procedure to less complex clinical laboratories that have the  
462 technology.

463 In some cases, a difference in the presence of a peak was observed in the  
464 duplicate of the same sample, which may be due to operator errors or by using  
465 the direct method, which presents greater variability than the extraction  
466 techniques. Evidence from the literature suggests that the protein extraction  
467 method extends or improves the range of peaks identified [16,20]. However, in  
468 this work direct method was prioritized, because it is much simpler, faster and  
469 easily applicable in the routine of any clinical laboratory. Due to the variability of



470 the method evidenced on some challenged spectra, the importance of working  
471 with technical replicates and analysing the presence/absence of several  
472 characteristic peaks is highlighted, in order to reduce the risk of making errors  
473 during the classification of one class or another. It is also evident that the  
474 search for a single peak is not enough, but rather the joint analysis of a profile,  
475 either manually or automatically, in order to performed a reliable identification.  
476 According to previous works by Mazzeo, 2006 and Teruyo, 2014, it was  
477 possible to confirm that the absence of the peak at 9060 Da is a useful indicator  
478 of STEC O157:H7, although it is not definitive per se, as could be evidenced in  
479 our work, where the finding of one or more of the nine detected peaks would  
480 confirm the presence of a STEC O157:H7 isolate, because these are not  
481 normally found in the other diarrheal types of the genus.  
482 The total absence of biomarker peaks that occurred in 4 isolates (three of them  
483 NON STEC O157 that would generally present a single peak) may also be due  
484 to errors in the technical procedure, which, in these cases, would be convenient  
485 to repeat or confirm by other methods if only this approach is available.  
486 A particular case occurred where the presence of 4 peaks was detected, as  
487 occurs on isolate of STEC O157, but with the presence, in addition to the  
488 representative peak of the DEC group at m/z 9060 Da, this isolate  
489 corresponded to a O157 H7 non-toxigenic type, which had been isolated from a  
490 meat processing plant, this strengthens the idea of being in the presence of a  
491 STEC isolate O157:H7 with the loss of the Shiga toxin phage, in the same way  
492 as described in other works **[21-22]**. Said isolate also presented other virulence  
493 factors such as enterohemolysin and eae, typical of STEC O157:H7.

494 Finally, the approach based on the detection of the presence/absence of peaks,  
495 although it is a manual method that requires a longer analysis time, presented  
496 excellent performance values and there were no differences regarding in  
497 sensitivity and specificity compared to the mathematical model, added to the  
498 availability of the Flex Analysis software in the equipment, the detection of these  
499 biomarker peaks could be applied in laboratories as a rapid screening method  
500 for suspicious STEC O157:H7 isolates.

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521 **REFERENCES.**

522

523 1. Rivas, M., Miliwebsky, E., Chinen, I., Deza, N., & Leotta, G. A. (2006).

524 Epidemiología del síndrome urémico hemolítico en Argentina.

525 Diagnóstico del agente etiológico, reservorios y vías de transmisión.

526 Medicina, 66(supplement 3), 27-32.

527 2. Miliwebsky E, Schelotto F, Varela G, Luz D, Chinen I, Piazza RMF.

528 Human Diarrheal Infections: Diagnosis of Diarrheagenic Escherichia coli

529 Pathotypes Escherichia coli in the Americas. Torres AG (ed.). En: Escherichia

530 coli in the Americas. Springer E-Book. ISBN: 978-3-319-45092-6. 2016.

531 Chapter 15. p. 343-69.

532 3. Kallow, W., Erhard, M., Shah, H.N., Raptakis, E. and Welker, M. (2010).

533 MALDI-TOF MS for Microbial Identification: Years of Experimental

534 Development to an Established Protocol. In Mass Spectrometry for

535 Microbial Proteomics (eds H.N. Shah and S.E. Gharbia).

536 <https://doi.org/10.1002/9780470665497.ch12>

537 4. He Y, Li H, Lu X, Stratton CW, Tang YW. Mass spectrometry biotyper

538 system identifies enteric bacterial pathogens directly from colonies grown

539 on selective stool culture media. J Clin Microbiol. 2010;48(11):3888-

540 3892. doi:10.1128/JCM.01290-10

- 541 5. Cherkaoui A, Hibbs J, Emonet S, et al. Comparison of two matrix-  
542 assisted laser desorption ionization-time of flight mass spectrometry  
543 methods with conventional phenotypic identification for routine  
544 identification of bacteria to the species level. *J Clin Microbiol.*  
545 2010;48(4):1169-1175. doi:10.1128/JCM.01881-09
- 546 6. Rocca MF, Zintgraff JC, Dattero ME, et al. A combined approach of  
547 MALDI-TOF mass spectrometry and multivariate analysis as a potential  
548 tool for the detection of SARS-CoV-2 virus in nasopharyngeal swabs. *J*  
549 *Virol Methods.* 2020;286:113991. doi:10.1016/j.jviromet.2020.113991
- 550 7. Miliwebsky E .Manual de procedimientos Escherichia coli productor de  
551 toxina Shiga en el marco de la detección de E. coli diarreigénico.2019.  
552 <http://sgc.anlis.gob.ar/handle/123456789/2307>
- 553 8. Espinal P, Seifert H, Dijkshoorn L, Vila J, Roca I. Rapid and accurate  
554 identification of genomic species from the *Acinetobacter baumannii* (Ab)  
555 group by MALDI-TOF MS. *Clin Microbiol Infect.* 2012;18(11):1097-1103.  
556 doi:10.1111/j.1469-0691.2011.03696.x
- 557 9. Bruker Daltonik GmbH, 2011.ClinPro Tools User Manual Version  
558 3.0.BrukerDaltonik GmbH, Bremen  
559 <https://doi.org/10.1371/journal.pone.0230334>.
- 560 10. Camoez, M., Sierra, J.M., Dominguez, M.A., Ferrer-Navarro, M., Vila, J.,  
561 Roca, I., 2016. Automated categorization of methicillin-resistant  
562 *Staphylococcus aureus* clinical isolates into different clonal complexes by  
563 MALDI-TOF mass spectrometry. *Clin.Microbiol. Infect.* 22 (2)  
564 <https://doi.org/10.1016/j.cmi.2015.10.009>, 161.e1-161. e7.

- 565 11. Zhang, H., Cao, J., Li, L., Liu, Y., Zhao, H., Li, N., Li, B., Zhang, A.,  
566 Huang, H., Chen, S., Dong, M., Yu, L., Zhang, J., Chen, L., 2015.  
567 Identification of urine protein biomarkers with the potential for early  
568 detection of lung cancer. *Sci.Rep.*5, 11805. <https://doi.org/10.1038/srep11805>.  
569
- 570 12. Wang, H.Y., Lien, F., Liu, T.P., Chen, C.H., Chen, C.J., Lu, J.J., 2018.  
571 Application of a MALDI-TOF analysis platform (ClinProTools) for rapid  
572 and preliminary report of sequence types in Taiwan. *PeerJ.* 6, e5784.  
573 <https://doi.org/10.7717/peerj.5784>
- 574 13. Stephens, M.A., 1974. EDF statistics for goodness of fit and some  
575 comparisons. *JASA* 69,730–737.
- 576 14. Landis, J., & Koch, G. The measurement of observer agreement for  
577 categorical data. *Biometrics*, 1977; 33, 159-174
- 578 15. Yumi Kubo, Osamu Ueda, Sawa Nagamitsu, Hachiro Yamanishi, Akihiro  
579 Nakamura, Masaru Komatsu, Novel strategy of rapid typing of Shiga  
580 toxin-producing *Escherichia coli* using MALDI Biotyper and ClinProTools  
581 analysis. *Journal of Infection and Chemotherapy*, Volume 27, Issue 8,  
582 2021, Pages 1137-1142, ISSN 1341-321X,  
583 <https://doi.org/10.1016/j.jiac.2021.03.002>.
- 584 16. Mazzeo Maria Fiorella, Alida Sorrentino, Marcello Gaita, Giuseppina  
585 Cacace, Michele Di Stasio, Angelo Facchiano, Giuseppe Comi, Antonio  
586 Malorni, and Rosa Anna Siciliano Matrix-Assisted Laser Desorption  
587 Ionization-Time of Flight Mass Spectrometry for the Discrimination of  
588 Food-Borne Microorganisms. *Appl Environ Microbiol.* 2006 Feb; 72(2):  
589 1180–1189. doi: 10.1128/AEM.72.2.1180-1189.2006.

- 590 17. Clark CG, Kruczkiewicz P, Guan C, et al. Evaluation of MALDI-TOF  
591 mass spectroscopy methods for determination of Escherichia coli  
592 pathotypes. J Microbiol Methods. 2013;94(3):180-191.  
593 doi:10.1016/j.mimet.2013.06.020
- 594 18. Fagerquist CK, Garbus BR, Miller WG, et al. Rapid identification of  
595 protein biomarkers of Escherichia coli O157:H7 by matrix-assisted laser  
596 desorption ionization-time-of-flight-time-of-flight mass spectrometry and  
597 top-down proteomics. Anal Chem. 2010;82(7):2717-2725.  
598 doi:10.1021/ac902455d
- 599 19. Ojima-Kato T, Yamamoto N, Suzuki M, Fukunaga T, Tamura H.  
600 Discrimination of Escherichia coli O157, O26 and O111 from other  
601 serovars by MALDI-TOF MS based on the S10-GERMS method. PLoS  
602 One. 2014;9(11):e113458. Published 2014 Nov 20.  
603 doi:10.1371/journal.pone.0113458
- 604 20. Ochoa ML, Harrington PB. Immunomagnetic isolation of  
605 enterohemorrhagic Escherichia coli O157:H7 from ground beef and  
606 identification by matrix-assisted laser desorption/ionization time-of-flight  
607 mass spectrometry and database searches. Anal Chem.  
608 2005;77(16):5258-5267. doi:10.1021/ac0502596
- 609 21. Bielaszewska M, Köck R, Friedrich AW, et al. Shiga toxin-mediated  
610 hemolytic uremic syndrome: time to change the diagnostic paradigm?.  
611 PLoS One. 2007;2(10):e1024. Published 2007 Oct 10.  
612 doi:10.1371/journal.pone.0001024
- 613 22. Friedrich AW, Zhang W, Bielaszewska M, et al. Prevalence, virulence  
614 profiles, and clinical significance of Shiga toxin-negative variants of

615           enterohemorrhagic Escherichia coli O157 infection in humans. Clin Infect

616           Dis. 2007;45(1):39-45. doi:10.1086/518573

617

618

619

620

621

622

623

624

625

626

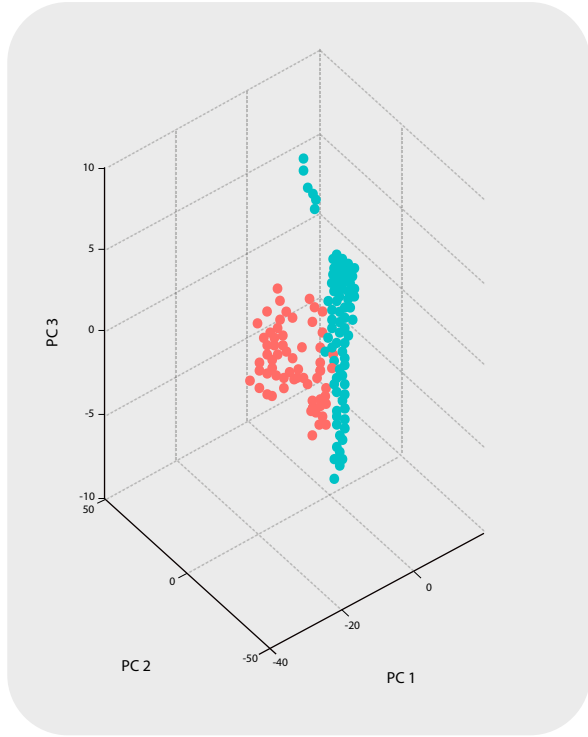
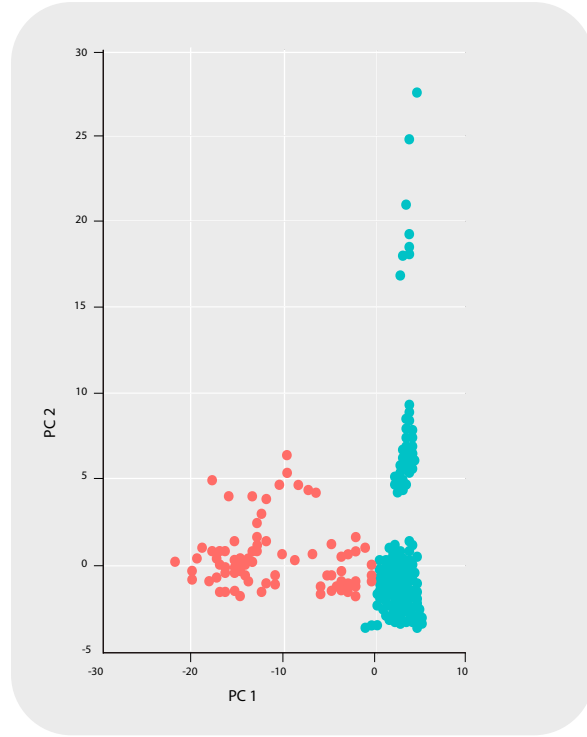
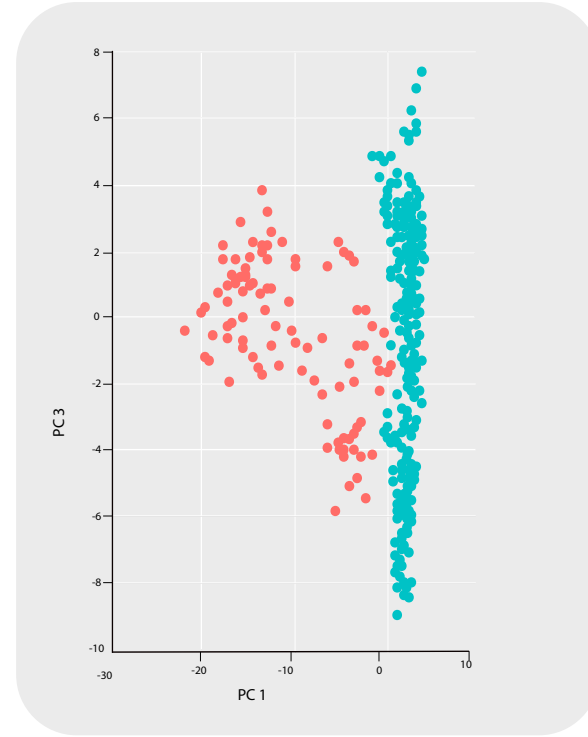
627

628

629

630

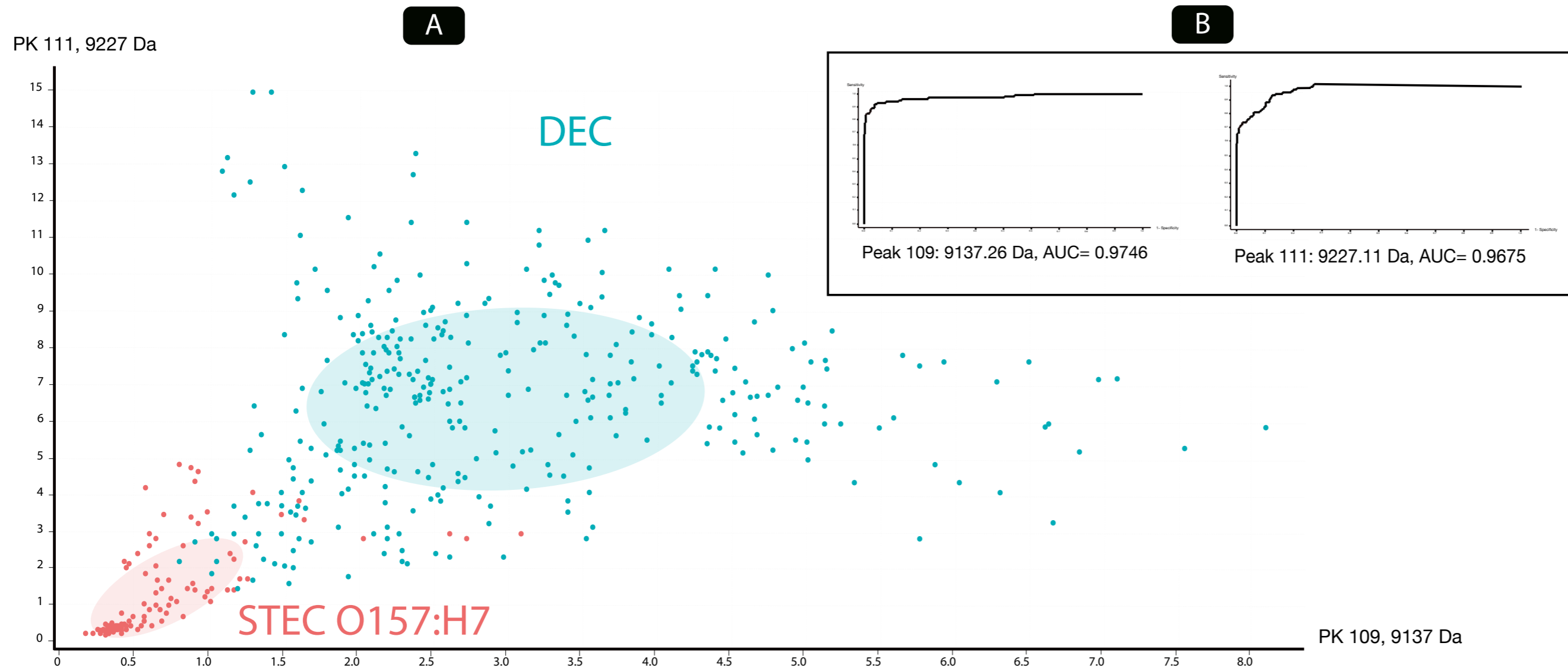
631

**A****B****C**

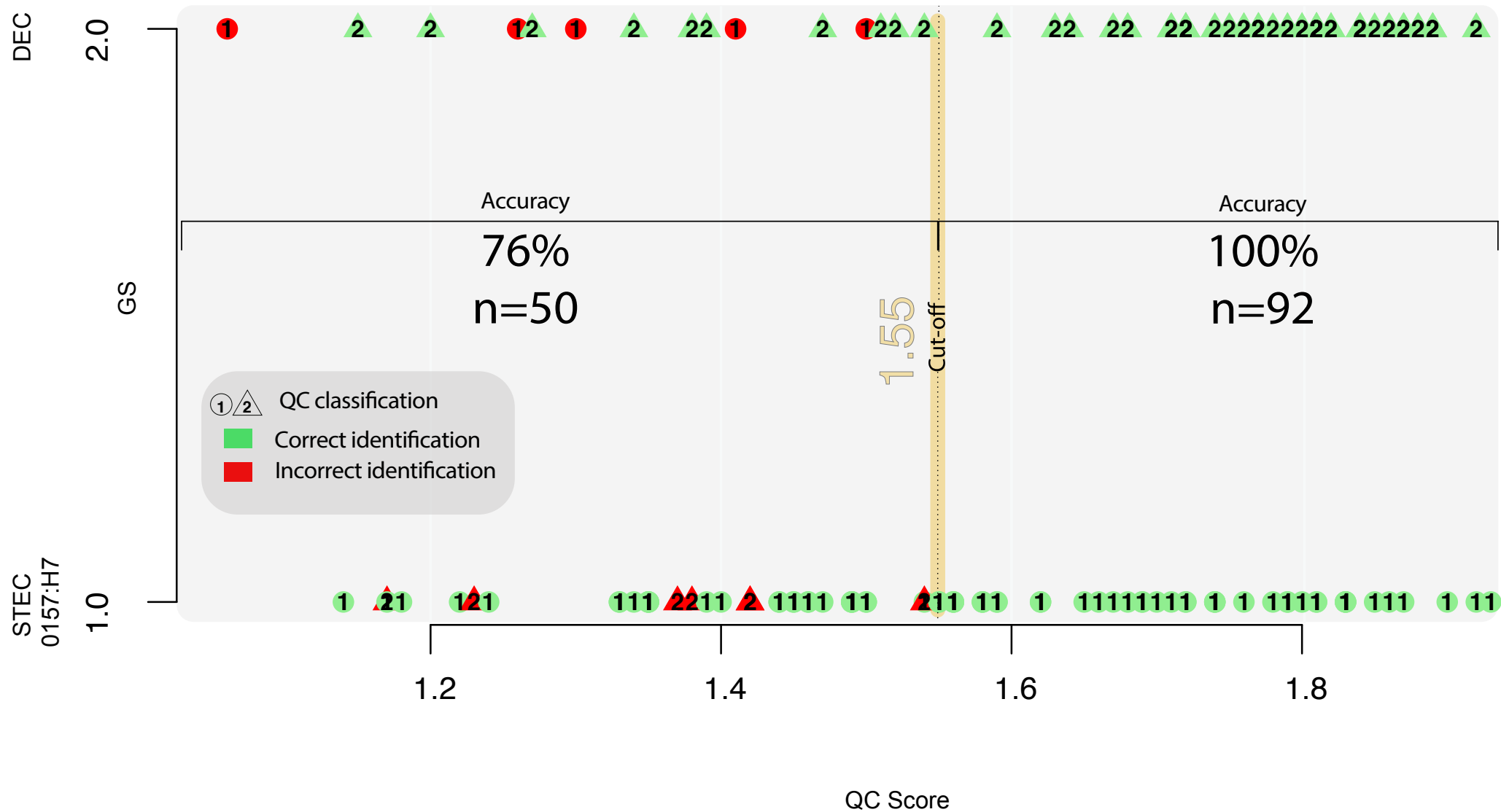
·DEC  
·STEC O157:H7

**Figure 1.** PCA plots results. Data originated from the external MATLAB software tool integrated into ClinPro Tools.

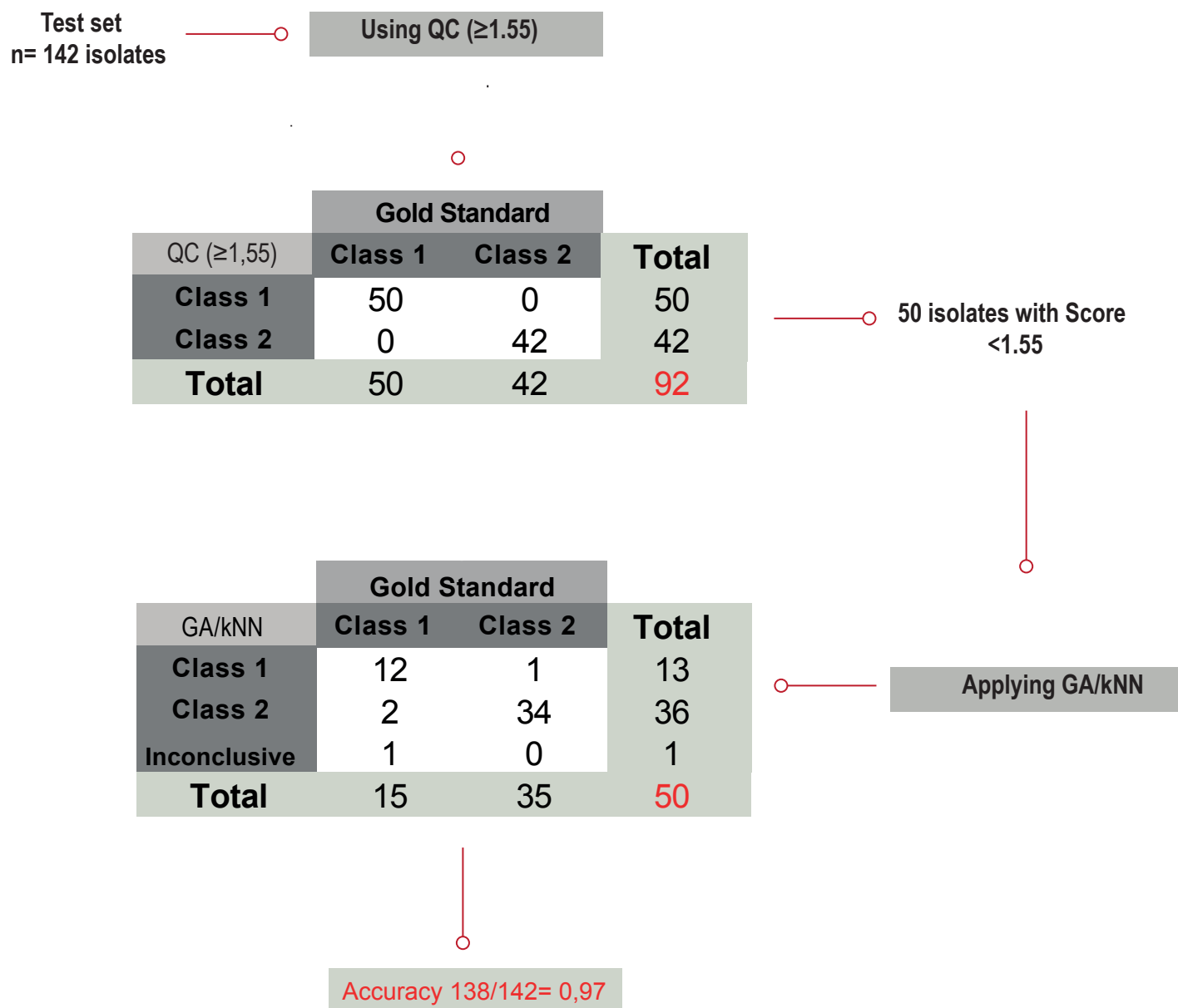




**Figure 2. A-** 2D graph of peak distribution of the 2-class model. **B-** ROC curves of the most discriminating peaks according to the analyses performed.

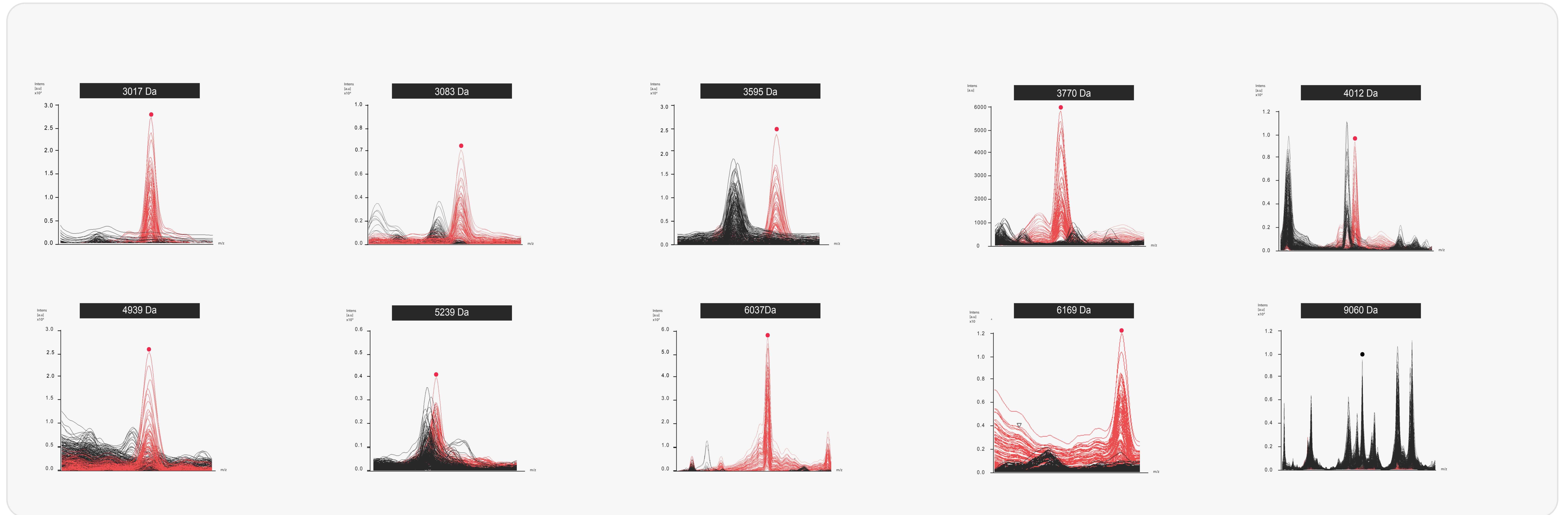


**Figure 3.** Graph of QC scores where the 100% concordance of the QC algorithm with the reference method is observed from the standardized cut-off value of 1.55.

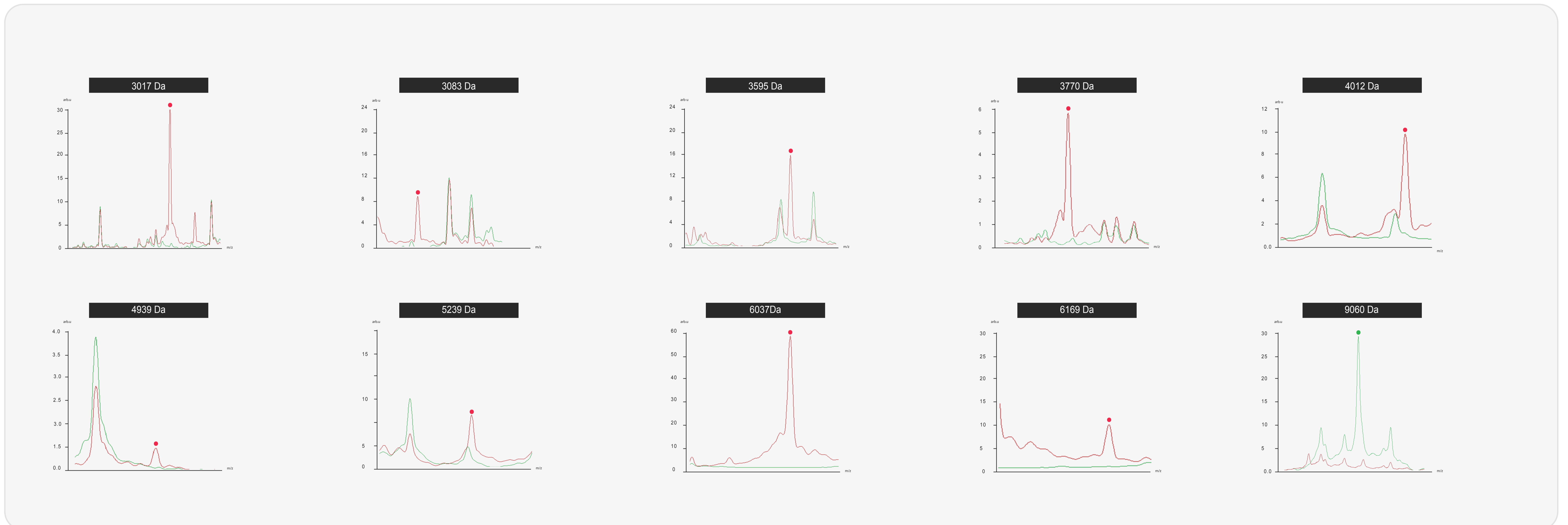


**Figure 4.** Algorithm applied for the identification of STEC O157:H7 based on predictive classification models.

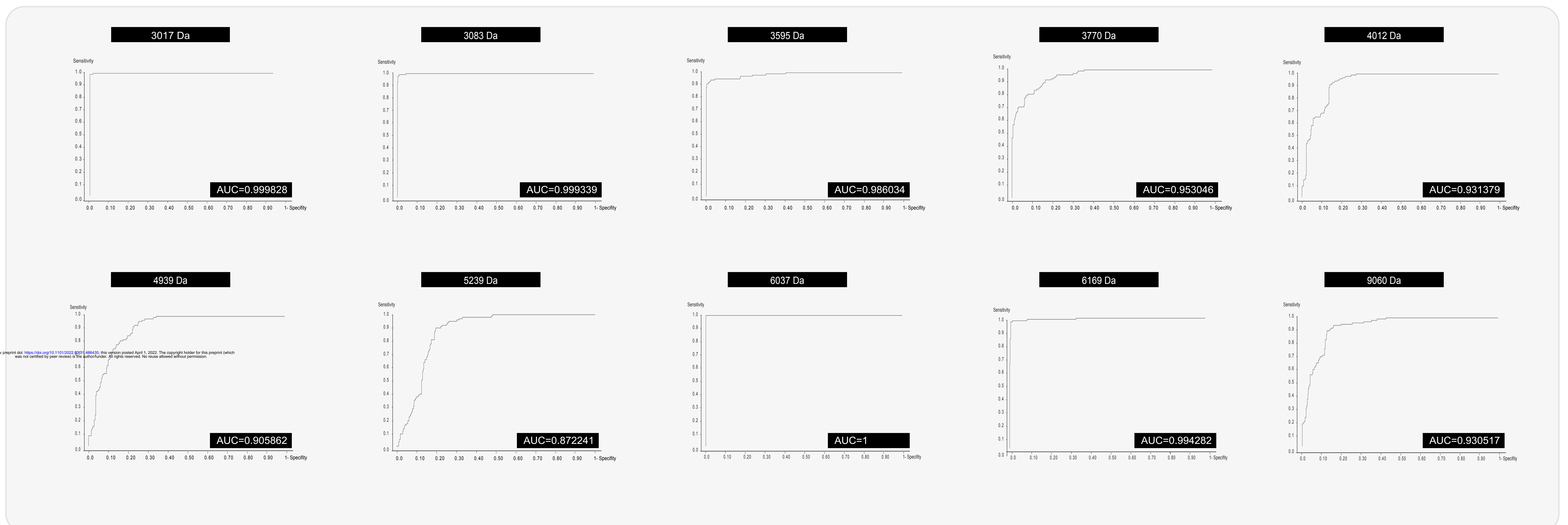
# A



# B



# C



**Figure 5.** A- Characteristic peaks (biomarkers) in individual spectra of STEC O157:H7 samples (red) versus DEC samples (black), obtained by manual analysis in Flex Analysis v3.4 software. B- Average spectra of the same peaks, STEC O157:H7 (red); DEC (green), obtained by ClinPro Tools v3.0. C- ROC curves and AUC values originated from the external MATLAB software tool integrated in ClinPro Tools.