

Anal Bioanal Chem (2013) 405:8363–8375
DOI 10.1007/s00216-013-7206-5

Comparison of Multianalyte Proficiency Test Results by Sum of Ranking 1 Differences, Principal Component- and Hierarchical Cluster Analysis

Biljana Škrbić, Károly Héberger, Nataša Đurišić-Mladenović

1 **Comparison of Multianalyte Proficiency Test Results by Sum of Ranking Differences,**
2 **Principal Component- and Hierarchical Cluster Analysis**

3 Biljana Škrbić^a, Károly Héberger^{b,*}, Nataša Đurišić-Mladenović^a

4 ^a University of Novi Sad, Faculty of Technology, Bulevar cara Lazara 1, 21000 Novi Sad, Serbia

5 ^b Research Centre for Natural Sciences, Hungarian Academy of Sciences, H-1025 Budapest,
6 Pustaszeri út 59-67, Hungary

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25 * Corresponding author: Károly Héberger, E-mail: heberger.karoly@ttk.mta.hu

26

27 **ABSTRACT**

28 Sum of ranking differences (SRD) was applied for comparing multianalyte results obtained by
29 several analytical methods used in one or in different laboratories, i.e. for ranking the overall
30 performances of the methods (or laboratories) in simultaneous determination of the same set of
31 analytes. The data sets for testing of the SRD applicability contained the results reported during one
32 of proficiency tests (PTs) organized by EU Reference Laboratory for Polycyclic Aromatic
33 Hydrocarbons (EU-RL-PAH). In this way, the SRD was also tested as a discriminant method
34 alternative to existing average performance scores used to compare multianalyte PT results. SRD
35 should be used along with the z -scores – the most commonly used PT performance statistics.
36 SRD was further developed to handle the same rankings (ties) among laboratories. Two benchmark
37 concentration series were selected as reference: (i) the assigned PAH concentrations (determined
38 precisely beforehand by the EU-RL-PAH), (ii) the averages of all individual PAH concentrations
39 determined by each laboratory.

40 Ranking relative to the assigned values and also to the average (or median) values pointed the same
41 laboratories with the most extreme results, as well as revealed groups of laboratories with similar
42 overall performances. SRD reveals differences between methods or laboratories even if classical
43 test(s) cannot. The ranking was validated using comparison of ranks by random numbers (a
44 randomization test) and using seven folds cross-validation, which highlighted the similarities among
45 the (methods used in) laboratories. Principal component analysis and hierarchical cluster analysis
46 justified the findings based on SRD ranking/grouping. If the PAH-concentrations are row-scaled
47 (i.e. z -scores are analyzed as input for ranking) SRD can still be used for checking the normality of
48 errors. Moreover, cross-validation of SRD on z -scores groups the laboratories similarly. The SRD
49 technique is general in nature, i.e. it can be applied to any experimental problem in which the
50 multianalyte results obtained either by several analytical procedures, analysts, instruments, or
51 laboratories need to be compared.

52 **Keywords:** ~~interlaboratory~~ multianalyte results, comparison, ~~polycyclic aromatic hydrocarbons,~~
53 ~~round-robin test,~~ sum of ranking differences, principal component analysis, comparison of ranks by
54 random numbers
55

56 **1. Introduction**

57

58 Due to public health concerns there has been a need in different domains (e.g. food safety,
59 environmental protection) for development of analysis that can identify and measure the numerous
60 contaminants belonging to the same or similar chemical groups in order to get as many as possible
61 data in one analytical run for the risk assessment. For instance, there are several lists of
62 contaminants belonging to different chemical classes (polycyclic aromatic hydrocarbons (PAHs),
63 polychlorinated biphenyls (PCBs), organochlorine pesticides, etc.) required or advised to be
64 monitored in food and/or environmental samples. In response to this, a number of methods have
65 been developed and applied routinely for the control of contaminants levels. Those methods that
66 can identify and measure a number of analytes concurrently are called “multianalyte (i.e.
67 multiresidue) methods” [1]. Although these methods are in routine use, they are often quite complex
68 and differ among themselves in terms of the sample preparation step, instrumental techniques
69 available, applied working parameters, etc. Multianalyte methods require not only careful
70 performance but also continuous monitoring to check the reliability of the measurements [1].
71 In order to verify the confidence in measurement results (or the competence of the laboratory either
72 accredited or non-accredited), including such multianalyte results, there is a request for the
73 laboratories to have quality control procedures for monitoring the performances of the analysis
74 undertaken (ISO/IEC 17025). One of the means to monitor the laboratory performance is its
75 participation in interlaboratory comparison programs. In an interlaboratory comparison ~~experiment~~,
76 different laboratories determine some characteristic, e.g. concentration of the same analyte(-s) in
77 one or various homogenous samples under documented conditions, assuming that the systematic
78 errors of methods in different laboratories follow normal distribution [2-6]. For simplicity, we use
79 the term interlaboratory comparison further on knowing that it is essentially analytical methods
80 comparison, (c.f. Table 1). The typical purposes for interlaboratory comparisons include evaluation
81 of the performance of laboratories for specific measurements, identification of problems in

82 laboratories and initiation of actions for improvement, establishment of the effectiveness and
83 comparability of test or measurement methods, provision of additional confidence to laboratory
84 customer, etc. [6]. In general there are two sub-types of interlaboratory studies: i) collaborative
85 trials or method performance (used to check the performance, generally the precision) of a single
86 analytical method, and ii) proficiency testing or laboratory performance studies (sometimes, the
87 term “round robin test” is also used) [3].

88 The laboratories participating in proficiency tests receive test material from the proficiency testing
89 provider; the material should be analyzed by measurement procedure of the choice, which is
90 consistent with the routine procedure in the laboratory. In the specified time period, the results of
91 the test material analysis should be reported to the proficiency testing provider, who further analyze
92 the results by appropriate statistical methods, generating summary statistics and performance
93 statistics in order to aid interpretation and to allow comparison with defined objectives. In fact, the
94 purpose is to measure deviation from the assigned value – a value attributed to a particular property
95 of a proficiency test material (e.g. concentration of analyte(-s)). Determination of the assigned
96 values belongs to the responsibility of the proficiency testing providers. The assigned value is not
97 disclosed to the participants until they have reported their results. Different statistical methods may
98 be used for calculation of the performance statistics; generally simple numerical or graphical
99 criteria, described in ISO 13528 [5] and ISO/IEC17043 [6] have been used to interpret the results
100 reported by laboratories participating in a proficiency test. The majority of these performance
101 statistics are generated from the results referring to the single analyte. If several analytes are
102 subjects of the proficiency test, performance statistics are generally given for each analyte
103 separately (i.e. the results for each analyte are analyzed separately). Additionally, in the case of
104 results for several analytes in one proficiency test material (multianalyte results), the use of some
105 graphical methods are recommended by ISO 13528 [5], describing the conditions and limitations of
106 these approaches. Youden [2] also describes a protocol how to complete an interlaboratory
107 examination, how to present data and what to do with the problems arisen (missing data, outliers,

108 and ranking laboratories). Youden suggests an approximate test to decide whether a laboratory
109 “passed” the test in measuring a single analyte (i.e. if produces acceptable results or not). The test is
110 based on *sum of ranks* and a special table with *critical values* could be used for comparison only in
111 the case if the number of participating laboratories is 15 or less. The other limitation of Youden
112 protocol is the number of objects (e.g. compounds content), which are also restricted and decision
113 on the laboratory accuracy needs a more sophisticated evaluation.

114 One of the most commonly used performance statistics is the z -scores calculated by Equation (1):

$$115 \quad z = \frac{(x-X)}{s} \quad (1)$$

116 where x is the participant’s result, X is the assigned value and s is the sample standard deviation for
117 proficiency assessment, which can be calculated by applying one of five proposed approaches [5,6].
118 The standard deviation for proficiency testing is used to assess laboratory bias, i.e. deviation from
119 the assigned value found in a proficiency test [5].

120 “Satisfactory” performance is indicated by an absolute value of z -score less or equal to 2. Absolute
121 values of z -score between 2 and 3 suggest “questionable” performance, while results are considered
122 “unsatisfactory” if absolute values of z -scores are above 3.0.

123 However, some authors highlighted that the z -score statistics can present pitfalls and have
124 limitations, so they should be interpreted cautiously [7,8].

125 Organization of the interlaboratory comparisons (ILCs) (PAHs) in food is one of the core duties of
126 the European Union Reference Laboratory for PAHs in food (EU-RL-PAH) hosted at the Institute
127 for Reference Materials and Measurements (IRMM) of the European Commission’s Joint Research
128 Centre. PAHs are a group of about ten thousand compounds, a few of them occurring in
129 considerable amounts in the environment and food, many being classified as probable or possible
130 human carcinogens. Human beings are exposed to PAHs mostly by intake of food, which is also the
131 reason why reliable analysis of PAHs in foodstuffs is of great importance. The activities of EU-RL-
132 PAH refer to Commission Regulation (EC) 1881/2006 [9] as amended by Commission Regulation
133 835/2011 [10] setting maximum levels of selected PAHs in various types of food, and to

134 Commission Regulation 333/2007 [11] as amended by Commission Regulation 836/2011 [12]
135 laying down sampling and analysis measures for the official control of the selected PAH levels in
136 foodstuffs.

137 Till now, there have been nine rounds of ILCs organized by EU-RL-PAH for 15+1 EU priority
138 PAHs (5-Methylchrysene-5MC, Benzo[a]anthracene-BAA, Benzo[a]pyrene-BAP,
139 Benzo[b]fluoranthene-BBF, Benzo[c]fluorine-BCF, Benzo[ghi]perylene-BGP,
140 Benzo[j]fluoranthene-BJF, Benzo[k]fluoranthene-BKF, Chrysene-CHR, Cyclopenta[cd]pyrene-
141 CPP, Dibenzo[a,e]pyrene-DEP, Dibenzo[a,h]anthracene-DHA, Dibenzo[a,h]pyrene-DHP,
142 Dibenzo[a,i]pyrene-DIP, Dibenzo[a,l]pyrene-DLP, and Indeno[1,2,3-cd]pyrene-ICP) in various
143 matrices, e.g. olive oils, sausages, solvent solutions, etc. Reports of these ILCs are readily available
144 on the official web site of IRMM:
145 http://irmm.jrc.ec.europa.eu/interlaboratory_comparisons/Pages/index.aspx.

146 These ILC studies aimed to evaluate trueness and precision of analytical results reported by the
147 participating laboratories for compounds belonging to the group of 15+1 EU priority PAHs in
148 different food matrices and to assess the influence of standard preparation and instrument
149 calibration on the performance of individual laboratories. The ILCs organized by EU-RL-PAH till
150 2010 have been designed and evaluated along the guidelines given in well approved ISO/IEC Guide
151 43 [13], while the latest proficiency tests have been conducted in accordance with ISO/IEC 17043
152 [6]. Additionally, the IUPAC International Harmonized Protocol for the Proficiency Testing of
153 Analytical Chemistry Laboratories has been also used in all proficiency tests of EU-RL-PAH [14].
154 The performance of the laboratories in determination of the target PAHs in selected food items
155 during the proficiency tests organized by EU-RL-PAH has been evaluated by z-score (Eq.1), in
156 which standard deviation for proficiency testing, s_p , for benzo[a]pyrene has been set to be equal to
157 the maximum tolerated standard measurement uncertainty, U_f , as defined by Commission
158 Regulation (EC) No 333/2007 [11] amended by Regulation (EC) 836/2011 [12]:

159
$$U_f = ((LOD/2)^2 + (\alpha C)^2)^{0.5} \quad (2)$$

160 where LOD relates to the required limit of detection (which is $0.3 \mu\text{g kg}^{-1}$ [11,12]), α is a numeric
161 factor depending on the concentration C (for C less or equal to $50 \mu\text{g kg}^{-1}$, α is 0.2 [11,12]). For
162 instance, the application of Eq. 2 with the assigned value of $3.0 \mu\text{g kg}^{-1}$ for benzo[a]pyrene and the
163 required limit of detection of $0.3 \mu\text{g kg}^{-1}$ results in a U_f value of $0.62 \mu\text{g kg}^{-1}$ (i.e. 20.6% of the
164 assigned value of $3.0 \mu\text{g kg}^{-1}$). For all other PAHs in the group of 15+1 EU priority PAH
165 compounds, standard deviation for proficiency testing was set to 22% of the assigned values of the
166 compounds of interest, as suggested by Thompson [15].

167 In this way, z -scores obtained for each analyzed PAH was used to assess the performance of the
168 laboratory (i.e. analytical method) taking into account PAH-compounds separately. Usually, bar-
169 plots of the z -scores grouped for each participating laboratory, have been used for visualisation of
170 the overall performance of the laboratories to analyze simultaneously all 16 PAHs. Such bar-plots
171 reveal common features in the z -scores for a laboratory (for instance, if one laboratory achieved
172 several high z -scores (higher than 2), a bar-plot would easily indicate a laboratory with poor
173 performance for these analyzed PAH compounds) [5]. Besides bar-plots, ISO 13528 [5] and
174 ISO/IEC 17043 [6] recommend use of other graphical methods in case of multianalyte proficiency
175 testing results, which combine performance scores for all analytes. For example, histogram type
176 plot of z -scores is a suitable method, when the number of measured characteristics is small. An
177 individual participating laboratory is identified by the position of its scores, which are used to assess
178 the lab performance. Nevertheless, these two documents discourage application of composite or
179 averaged performance scores (e.g. average absolute z -score) because they can mask poor
180 performance on one or more analytes, also suggesting that simply the number (or percentage) of
181 results determined to be acceptable could be used in case of multianalyte proficiency tests.

182 There has also been an attempt to improve well established combined z -scores for evaluation of the
183 overall laboratory performance in application of multianalyte method [8]. There is a definite
184 scarcity of the works on introducing alternatives to the existing procedures for assessment of the
185 laboratory performance in multianalyte determination. Thus, the aim of this work is to contribute to

186 those scarce alternatives and to test a simple ~~alternative~~ method based on sum of ranking
187 differences (SRD) for comparative assessment of the overall performances of laboratories in
188 multianalyte determinations.

189 SRD is simple, entirely general technique suitable to order methods, models, to find their
190 similarities and the differences among them [16]. The SRD procedure is easy to apply and it
191 provides a unique ranking. So far, this technique (SRD) has been used in different fields (e.g. for
192 column selection in chromatography [16], for selection of the best polarity measure for small
193 organic molecules [17], for sensory panel testing [18-20], for comparison and ranking of
194 QSAR/QSPR models, including selection of metric for QSAR models [21-24], for PLS model
195 comparison in near infra-red spectroscopy [25], for testing performance for Raman spectra
196 resolution [26], for Hansen's solubility parameters [27], for comparison of biochemical assay (Elisa
197 veratox) and liquid chromatography in determination of mycotoxin contents [28], for comparative
198 evaluation of acidic dissociation constants [29]. There has not been any attempt to apply it for
199 comparison of analytical results obtained in different laboratories, including also those from
200 interlaboratory comparisons. Here we extend the SRD procedure to evaluate laboratories according
201 to the overall performance taking into account multianalyte results simultaneously not just
202 evaluating the quality in measuring one individual compound.

203 The data reported for 15+1 EU PAHs during the ILCs organized by EU-RL-PAH were taken for
204 testing this new technique; one of the major reasons for using these data is their availability and
205 abundance, providing the source for SRD validation on different data sets. In this way, the SRD was
206 also tested as a discriminant method alternative to existing average performance scores used to
207 compare multianalyte PT results. SRD should be used along with the z -scores, and it was compared
208 with well-known chemometric techniques, too. Additionally, the ranking was validated by
209 Comparison of Ranks of Random Numbers (CRNN procedure), which is a kind of permutation test
210 [16,30] and by leave-many-out cross-validation (CV) [31]. The ranking made by SRD was
211 compared to the results of principal component analysis and hierarchical cluster analysis.

212 2. Experimental

213 2.1. Data sets

214 Data published in Report on the 5th ILC for determination of 15+1 EU priority PAHs in edible oil
215 [32] were used to form the input matrices: 16 PAH-compounds (samples) were enumerated in the
216 rows, whereas laboratories (analytical methods) were arranged in columns and were coded as L1,
217 L2, ..., L13. The edible oil sample was provided by the ILC organizer and it was an olive oil spiked
218 with 15+1 EU priority PAHs. Of all laboratories ~~included~~ in the ILC, only those (13 laboratories)
219 that reported the results for all 16 PAHs of interest, were included in this study, since the input
220 matrices for SRD testing should be without empty cells, which is the case when results for some
221 PAH compounds were not reported. The laboratories were free in the selection of the test method
222 for sample preparation and PAH determination. The reported results, corresponding z-scores and
223 methods used by participating laboratories, taken from the report of the 5th ILC of EU-RL-PAH, are
224 summarized in Table 1. The percentages of acceptable results (z-scores less or equal to 2) are also
225 presented in Table 1 for each selected laboratory.

226 Table 1

227 Two data sets based on the experimental results [32] were formed for testing the applicability of
228 SRD procedure:

- 229 • “OIL” set was formed of the PAHs contents in edible oil sample reported by each
230 participating laboratory (“reported” results presented in Table 1); the set size was 16 rows
231 (PAH-compounds) \times 13 columns (laboratories or methods);
- 232 • “OIL+As” set was in fact the “OIL” set extended with the column containing the assigned
233 values – analyte concentrations in ILC test material (spiked edible oil sample) determined
234 beforehand by EU-RL-PAH (i.e. calculated from gravimetric preparation data); thus, its size
235 was 16×14 .

236 Furthermore, an additional data set, so-called “Z-SCORE” set (16×13), was created of the absolute
237 values of z-scores calculated by the ILC organizer using Eq. (1) (presented in Table 1 [32]).

238 2.2. *Sum of ranking differences*

239 The key step in SRD procedure is the selection of the reference for ranking, when the true (ideal,
240 benchmark) ranking is not known [16]. Often the ranking by average values can be accepted as
241 “ideal”, since the errors cancel each other. The maximum likelihood principle will ensure that the
242 most probable ranking will be provided by the average. The methods that deviate from the average
243 less are ranked ahead. The best ranking is not necessarily provided by the average values, as it can
244 be a known sequence (here the assigned values), the maximum (if comparing best classification
245 rates) or the minimum (in the case of error rates and residuals). For the sets created in this study the
246 following references for ranking of the laboratories ~~values~~ were chosen:

- 247 a) the assigned values of 15+1 EU PAHs contents in edible oil sample (last column in Table 1)
- 248 as a reference for ranking within “OIL” set,
- 249 b) the averages of the reported results (values in $\mu\text{g}/\text{kg}$ presented in Table 1) and the assigned
- 250 value for each compound (row averages) as a reference for “OIL+As” set, while
- 251 c) the minimums of the absolute values of z -scores for each compound (presented in Table 1;
- 252 row minimums) for “Z-SCORE” set.

253 These selection were the logical choices in order to test SRD procedure: a) ranking of the reported
254 values on the base of the known (assigned) values would indicate laboratories that obtained
255 multianalyte results most similar to the assigned values; b) similar indication might be expected if
256 the assigned values would be included into the input set and then using the “overall” averages of
257 reported and assigned results (which, by the way, could be assumed to converge to the true values),
258 leading to the simultaneous ranking of the assigned and reported values, and finally c) ranking of
259 the laboratories according to their absolute z -score values in comparison to the minimal (absolute)
260 z -scores (representing the minimum deviation from the assigned value) per each compound. The
261 absolute values of the z -scores would allow a direct estimation of the performance of the
262 laboratories, but calculation of z -scores realizes a row-standardization (c.f. Eq. (1)), i.e. differences
263 needed for ordering are destroyed by row standardization. Hence, the absolute values of z -scores

264 order the laboratories randomly, and hence they are suitable to check whether the initial assumption
265 (normality) is valid or not.

266 Each (individual) laboratory was ranked and compared to the above mentioned references in
267 following way: the ascending reference values of PAH concentrations were ordered giving them
268 consecutive numbers from 1 to 16 (this is so-called “reference (benchmark) ranking”). Then,
269 ranking of data within each column (i.e. ranking of the results of each laboratory) was made (so-
270 called “individual ranking”); the absolute values of the differences between the reference and the
271 individual rankings for all compounds were calculated and summed for each laboratory. In this way,
272 the sum of (absolute) ranking differences, SRD values, were calculated for each laboratory. The
273 closer is the SRD value to zero (i.e. the closer is the sum of differences of individual ranking to the
274 reference one) the better is the analytical method for simultaneous determination of all analytes.
275 The proximity of SRD values shows that the methods used by the laboratories have similar (overall)
276 performance in the multianalyte (PAHs) determination. Equal concentrations (so-called ties) to two
277 digits received the same rank number during the ranking procedure.

278

279 *2.3. Validation*

280 Two types of validations have been carried out (i) comparison of ranks by random numbers
281 (CRRN), which is in fact a randomization test [16,30], and (ii) leave-many-out (seven folds) CV
282 followed the literature recommendation [31]. Namely, (i) CRRN procedure includes the
283 determination of the theoretical distribution for ranking using solely random numbers and the
284 distribution is compared to the actual rankings; (ii) during the seven folds CV (approximately) 1/7
285 of the objects were left out and the ranking was made on the remaining 6/7th number of objects just
286 seven times. The different rankings provided uncertainties for the SRD values.

287

288 *2.4. Exploratory statistics*

289 In the exploratory phase box and whisker plots were used to graphically present numerical data like
290 z-scores and cross-validated SRD values, while hierarchical cluster analysis (HCA) and principal
291 component analysis (PCA) were applied on the above sets in order to observe the similarity and
292 dissimilarity of laboratories (methods), to analyze quantitatively the relationships among the results
293 of laboratories (i.e. their analytical efficiency) and to compare these results with the SRD ones.
294 Mean centering and scaled to unit standard deviation were applied as data preprocessing step before
295 principal component- and hierarchical cluster analysis. Standard procedures were applied
296 (StatisticaTM, version 7.0, StatSoft Inc., Tulsa Oklahoma, USA).

297

298 **3. Results and discussion**

299 *3.1. Exploratory statistics*

300 PCA shows (Figure 1) the grouping of the laboratories within the “OIL+As” set (thus, grouping
301 relative to—the assigned values similar grouping can be observed for absolute values of z-scores (“Z-
302 SCORES” set). Figures 1a and 1b show the loading plots of two main PCs retained in both cases
303 that accounted similar share of the total data variance (~70%). The L5 was by far the most outlying
304 laboratory when the reported values were compared to the assigned (Figure 1a); there were few
305 more points (L2, L6, L7, L8, L9, L12) diverging from the central cluster comprising of the
306 laboratories (L1, L3, L4, L10, L11, L13) closest to the assigned value. The score plot for the z-
307 scores (Figure 1b) also pointed out L5 as an outlier and similarities among L1, L3, L4, L10, L11,
308 and L13.

309

Figures 1a and 1b

310 The dendrogram of Figure 2 indicates clustering of the laboratories similarly to the PCA groupings.
311 Laboratories L6 and L9 and particularly L5 reported the most dissimilar results to those reported by
312 the other labs (Figures 2a and 2b) and also to the assigned values determined by the proficiency
313 testing provider, EU-RL-PAH (Figure 2a).

314

Figures 2a and 2b

315 The results reported by laboratories L1, L3, L4, L10, L11 and L13 form a dense cluster (the
316 assigned values also belong to this cluster on Figure 1a). The same pattern can be observed on the
317 PCA plots (Figures 1a and 1b).
318 Z-scores of these six laboratories were all below 2, while the rest of laboratories had at least one z-
319 score (its absolute value) higher than 2, indicating questionable ($2 < |z| \leq 3$) or unsatisfactory
320 ($|z| > 3$) performances for one (or more) particular PAH compound(s). Box and whisker plots of
321 the absolute values of z-scores of the laboratories are given in Figure 3. The outlying laboratory L5
322 could be easily seen in Figure 3a; after its exclusion (Figure 3b) the laboratories might be ordered
323 according to the median absolute values of z-scores as follows (median absolute values of z-scores
324 are given in parentheses): L4 (0.25) ~ L11 (0.265) < L3 (0.355) ~ L10 (0.37) < L1 (0.39) < L13
325 (0.455) < L9 (0.675) ~ L12 (0.685) < L2 (0.715) < L8 (0.885) < L7 (1.03) < L6 (2.41). Apart from
326 L5, the highest standard deviations (SD) of the absolute values of z-scores were observed for L6
327 and L9 (SD for both laboratories SD = 1.26), while for others, the SDs were in the range from 0.25
328 (for L13) to 0.78 (for L8).

329 Figures 3a and 3b

330

331 3.2. Sum of ranking differences

332 The SRDs calculated for “OIL” and “OIL+As” data sets can be seen in Figure 4Table 2. Similarities
333 (i.e. groupings) of laboratories can also be observed, as well as their dissimilarities from the
334 ordering point of view, i.e. SRD can also be considered as a dissimilarity measure (the higher its
335 value, the more dissimilar to the reference value) [16,30]. Thus, the best ranked laboratories
336 according to the lowest SRD values in “OIL” and “OIL+As” sets appeared to be L2 and L3 (Table
337 2Figure 4); they showed the best overall performance in simultaneous determination of 15+1 EU
338 PAHs.

339

Figure 4

340 It is interesting to note that both laboratories differed from others by using the sample preparation
341 method based on size-exclusion chromatography (gel permeation chromatography) followed by
342 high performance liquid chromatography with fluorescence detection (Table 1). It could also be
343 seen that proximity of the SRD values indicates similar performances in analyzing 15+1 EU priority
344 PAHs among majority of the laboratories (ten laboratories in “OIL” set had SRD between 3-11,
345 while in “OIL+As”, SRDs of eleven laboratories ranged from 8 to 14). Three laboratories (L9, L6
346 and L5) had distinguishable higher SRDs (Table 2 Figure 4) as a consequence of significantly lower
347 performances in analyzing 15+1 EU PAHs. The L5 was the worst ranked laboratory in “OIL” and
348 “OIL+As” sets; it should be noted that only this lab used method for determination of PAHs based
349 on liquid-liquid/solid-phase extraction followed by gas chromatography coupled to mass
350 spectrometry. The best two laboratories (L2, L3) are somewhat better than the assigned values (L10
351 is equivalent) if accepting the mean average value as reference for ranking within the “OIL+As” set.
352 The ranking of laboratories in these two sets, other than those ranked as “the best” and “the worst”,
353 was slightly different.

354 Even though L2 had one z -score (its absolute value) slightly higher than 2 ($z = 2.02$) it was ranked
355 exactly on the same way as L3, indicating that SRD procedure might conceal one result very close
356 to the questionable performance, but it clearly depicts the laboratories with the poorest
357 performances (outlier).

358 In order to check the influence of the outlier on the ranking in “OIL+As” set, the SRD procedure
359 was also applied on the set without L5 (so-called “OIL+As-OUT” set) and the resulting SRDs
360 (calculated on the base of the averages used for the reference ranking) are also presented in Table
361 2 Figure 4. The rationale behind this lies in fact that the average values selected for the reference
362 ranking in “OIL+As” set were directly affected by the all input values (including the outlier),
363 contrary to the reference chosen for “OIL” set (i.e. the assigned values cannot be influenced by the
364 presence of outlier). Removing the outlier (“OIL+As-OUT” set, using averages as the reference)
365 caused slightly less SRDs for L8 and L10 (Figure 4). An alternative would be the selection of

366 median (or other robust measure) instead of the averages triangles in Figure 4). Interestingly L13 is
367 ranked first (slightly better than the assigned values), which exhibits the smallest range on Figure
368 3b. Other patterns are mostly similar to the remaining rankings of-Figure 4. From the comparison of
369 SRD rankings Figure 4 it could be concluded that median is the best choice. Figure 4 also contains
370 the normalized sum of squared z -scores (SZ2norm, a Euclidean distance) suggested as the most
371 optimal overall performance indicator by Medina-Pastor et al. [8]. All indicators in Figure 4 were
372 placed on the same scale between 0 and 100. As a non-robust measure, SZ2norm is sensitive to the
373 outlying observation most. Almost all variability in the data (>94%) is carried by the L5 outlier.
374 Any variants of SRD ranking are robust and allowed observing differences in other laboratories as
375 well (on the expense of the outlying L5).
376 The SRDs for laboratories were scaled between 0 and 100 (Figures 5a and 5b) in order to be
377 comparable among each other [16]. It could easily be seen that the location of the scaled values for
378 majority of laboratories was far from the SRDs of random numbers in the case of “OIL” and
379 “OIL+As” sets (Figures 5a and 5b, respectively), showing that their ranking was far from being
380 random. The L5 was the worst ranked laboratory in “OIL” and “OIL+As” sets; its SRD value in
381 both sets was close to the first icosaille (5%).

382 Figures 5a and 5b

383 The SRDs calculated for the “Z-SCORE” set were quite different than those obtained for the
384 reported values (i.e. for “OIL”, “OIL+As”, and “OIL+As-OUT” sets), as expected, because row-
385 standardization eliminates the differences needed for ordering. However, the overlapping with
386 normal distribution for the z -scores can easily be seen on Figure 6. All SRD values for the “Z-
387 SCORES” set overlapped with random distribution, except for L3 (Figure 6), which was also
388 located very close to the first icosaille, indicating that ordering of labs based on the absolute z -scores
389 for all compounds is not better than the random ordering (ordering of random numbers). In order to
390 check this observation, the SRD with CRNN procedure was also used on absolute values of z -scores
391 calculated for 24 laboratories participating in the 7thILC on PAHs in edible oil [33] and for 14

392 laboratories during the 4th ILC on PAHs in fish [34], and, again, SRDs overlapped with the random
393 distribution (data not shown). These observations can be considered as a proof that the errors of labs
394 (i.e. the deviation of their results from the assigned values, not the individual PAH concentrations)
395 expressed as z -scores follow a normal distribution.

396 Figures 6a and 6b

397 To reveal uncertainties for SRD, cross-validation (seven-fold CV [30,31]) has been carried out. Box
398 and whisker plots clearly exhibit the difference between classical (statistical) and present (SRD)
399 approach (Figures 7a and 7b, respectively). Figure 7a allows observing one outlying laboratory (L5)
400 nothing else, whereas seven-fold CV of SRD values allow us to group the laboratories similarly to
401 Figures 2a, 2b and 5a.

402 Figures 7a and 7b

403 Figure 7b shows the same pattern as Figure 5a with subtle, negligible differences suggesting that
404 cross-validation does not change the ranking of laboratories just helps in grouping them.
405 Comparing the results of PCA, HCA, SRD and CV-SRD, shows the very same (or almost the same)
406 clustering pattern. Moreover, CV-SRD reveals the uncertainties in the ranking and clustering. Sign
407 test or Wilcoxon's matched pair test is suitable to decide about the significance of CV-SRD
408 grouping.

409

410 **4. Conclusions**

411 Sum of Ranking Differences methodology (SRD) is a simple technique general in nature that can be
412 ~~used as~~ applied to any experimental problem in which the multianalyte results obtained either by
413 several analytical procedures, analysts, instruments, or laboratories need to be compared. Besides
414 the z -scores, the most commonly used PT performance statistics that assess the results of each
415 analyte separately, SRD could be regarded as an alternative way for ranking of measurement
416 methods and laboratories involved in interlaboratory comparison tests according to their
417 multianalyte results. SRD provides similar groupings as classical techniques (principal component

418 and hierarchical cluster analysis) and it is more influential than the (normalized) sum of the squared
419 z-scores.

420 The overall bias covering simultaneously the results on the whole group of targeted analytes is
421 taken into account (the bias follows normal distribution). SRD takes the disadvantages of the earlier
422 evaluation methods out (e.g. the discrepancies in ranking for individual compounds).

423 SRD proved to be a useful tool in choosing the analytical methods or the laboratories with the best
424 overall performances in multianalyte determinations. An unambiguous selection of the
425 laboratory(ies) or analytical methods could be made that produce results the most similar to the
426 assigned values, if comparison of the overall (multianalyte) performances of laboratories
427 participating in PT programs is made. SRD could point out the method(s) that produce(s) the best
428 results with respect to the overall averages (or medians), if the comparison of several multianalyte
429 methods should be taken. Similarly, the laboratories with the most extreme results could be easily
430 pointed out in any of the above two cases. Additionally, grouping of laboratories with similar
431 overall (multianalyte) performances can be obtained in similar manner by multivariate techniques
432 such as principal component analysis and hierarchical cluster analysis.

433

434 **Acknowledgement**

435 The results presented here are the part of project no. 172050 "Development and application of the
436 advanced chromatographic and spectrometric methods in the analysis of xenobiotics and their
437 degradation pathways in biotic and abiotic matrices", coordinated by Prof. B. Škrbić and supported
438 by the Ministry of Education and Science of the Republic of Serbia, as well of the bilateral project
439 of the Hungarian-Serbian Intergovernmental S&T Cooperation Program for 2010–2011 "Comparison
440 of various analytical and chemometric methods" funded by the Ministry of Education and Science of
441 the Republic of Serbia and National Innovation Office of Hungary.

442

443

444 **References:**

- 445 [1] W. Horwitz, T. Jacksin, S.J. Chirtel, J. AOAC Int. 84 (2001) 919-935.
- 446 [2] W. J. Youden, Statistical Manual of the Association of Official Analytical Chemists: Statistical
447 Techniques for Collaborative Test, AOAC International, Gaithersburg, 1975.
- 448 [3] E. Hund, D. Luc Massart, J. Smeyers-Verbeke, Anal. Chim. Acta 423 (2000) 145–165.
- 449 [4] M. Thompson, R. Wood, J. AOAC Int. 76 (1993) 926–940.
- 450 [5] ISO 13528 – Statistical methods for use in proficiency testing by interlaboratory comparisons,
- 451 [6] ISO/IEC 17043 – Conformity assessment – General requirements for proficiency testing, 2010.
- 452 [7] M. Ricci, O. Bercaru, R. Morabito, C. Brunori, I. Ipolyi, C. Pellegrino, A. Sahuquillo, F.
453 Ulberth, TrAC - Trends Anal. Chem. 26 (2007) 818-827.
- 454 [8] P. Medina-Pastor, M. Mezcuca, C. Rodriguez-Torreblanca, A.R. Fernandez-Alba, Anal. Bioanal.
455 Chem. 397 (2010) 3061-3070.
- 456 [9] COMMISSION REGULATION (EC) No 1881/2006 of 19 December 2006 setting maximum
457 levels for certain contaminants in foodstuffs, Official Journal of the European Union 364,
458 20.12.2006.
- 459 [10] COMMISSION REGULATION (EU) No 835/2011 of 19 August 2011 amending Regulation
460 (EC) No 1881/2006 as regards maximum levels for polycyclic aromatic hydrocarbons in
461 foodstuffs, Official Journal of the European Union L 215, 20.8.2011.
- 462 [11] COMMISSION REGULATION (EC) No 333/2007 of 28 March 2007 laying down the
463 methods of sampling and analysis for the official control of the levels of lead, cadmium,
464 mercury, inorganic tin, 3-MCPD and benzo(a)pyrene in foodstuffs. Official Journal of the
465 European Union L 88, 29.3.2007.
- 466 [12] COMMISSION REGULATION (EU) No 836/2011 of 19 August 2011 amending Regulation
467 (EC) No 333/2007 laying down the methods of sampling and analysis for the official control
468 of the levels of lead, cadmium, mercury, inorganic tin, 3-MCPD and benzo(a)pyrene in
469 foodstuffs, Official Journal of the European Union L 215, 20.8.2011

- 470 [13] ISO/IEC Guide 43-1:1997 (E), Proficiency testing by interlaboratory comparisons: Part 1:
471 Development and operation of proficiency testing schemes
- 472 [14] M. Thompson, S.L.R. Ellison, R. Wood, *Pure Appl. Chem.* 78 (2006) 145–196.
- 473 [15] M. Thompson, *Analyst* 125 (2000) 385–386.
- 474 [16] K. Héberger, *TrAC - Trends Anal Chem*, 29(2010)101–109.
- 475 [17] K. Héberger, I.G. Zenkevich, *J. Chromatogr. A*, 1217 (2010) 2895–2902.
- 476 [18] K. Kollár-Hunek, J. Heszberger, Z. Kókai, M. Láng-Lázi, E. Papp, *J. Chemometr.* 22 (2008)
477 218–226.
- 478 [19] L. Sipos, Z.Kovács, D. Szöllösi, Z. Kókai, I. Dalmádi; A. Fekete, *J. Chemometr.* 25 (2011) 275–
479 286.
- 480 [20] V. Losó, A. Tóth, A. Gere, J. Heszberger, G. Székely, Z. Kókai, L. Sipos, *Acta Alimentaria* 41
481 (2012) 109–119 .
- 482 [21] Z. Garkani-Nejad, M. Ahmadvand, *Chromatographia* 73 (2011) 733–742.
- 483 [22] X. Liu, Y. Ren, P. Zhou, Z. Shang, *J. Mol. Struct.* 995 (2011) 163–172.
- 484 [23] K. Héberger, B. Škrbić, *Anal. Chim. Acta* 716 (2012) 92–100.
- 485 [24] K. Roy, I. Mitra, P. K. Ojha, S. Kar, R. N. Das, H. Kabir, *Chemometr. Intell. Lab. Syst.* 118
486 (2012) 200–210.
- 487 [25] A. A. Gowen, G. Downey, C. Esquerre, C. P. O'Donnell, *J. Chemometr.* 25 (2011) 375–381.
- 488 [26] B. Vajna, A. Farkas, H. Pataki, Z. Zsigmond, T. Igricz, G. Marosi, *Anal. Chim. Acta*, 712
489 (2012) 45–55.
- 490 [27] K. Bielicka-Daszkiwicz, A. Voelkel, K. Héberger, M. Pietrzyńska, *J. Chromatogr. A* 1217
491 (2010) 5564–5570.
- 492 [28] E. K. Tangni, J.–C. Motte, A. Callebaut, A. Chandelier, M. De Schrijver, L. Marnix, L.
493 Pussemier, *Mycotoxin Res.* 27 (2011) 105–113.
- 494 [29] G. T. Balogh, A. Tarcsay, G. M. Keserű, *J. Pharm. Biomed. Anal.*, 67–68 (2012) 63–70.
- 495 [30] K. Héberger, K. Kollár-Hunek, *J. Chemometr.* 25 (2011) 151–158.

- 496 [31] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning: Data Mining,
497 Inference, and Prediction, 2nd ed., Springer, New York, 2009.
- 498 [32] D. Lerda, L. Hollósi, P. Lopez, S. Szilágyi, T. Wenzl, Report on the 5th interlaboratory
499 comparison test organized by the European Union Reference Laboratory for Polycyclic
500 Aromatic Hydrocarbons-15+1 EU priority PAHs in edible oil and acetonitrile, JRC Scientific
501 and Technical Reports, European Commission, Joint Research Centre, Institute for Reference
502 Materials and Measurements, 2010.
- 503 [33] D. Lerda, P. L. Sanchez, S. Szilágyi, P. Verlinde, T. Wenzl, Report on the 7th inter-laboratory
504 comparison test organised by the European Union Reference Laboratory for Polycyclic
505 Aromatic Hydrocarbons “15+1 EU priority PAHs in spiked olive oil and solvent
506 solution”,JRC Scientific and Technical Reports, European Commission, Joint Research
507 Centre, Institute for Reference Materials and Measurements,2011.
- 508 [34] D. Lerda, L. Hollósi, P. Lopez, S. Szilágyi, T. Wenzl, Report on the 4th interlaboratory
509 comparison test organised by the Community Reference Laboratory for Polycyclic Aromatic
510 Hydrocarbons “15+1 EU Priority PAHs in fish and acetonitrile”,JRC Scientific and Technical
511 Reports, European Commission, Joint Research Centre, Institute for Reference Materials and
512 Measurements, 2010.
- 513 .

514 Table 1 Summary of data (reported results in µg/kg, z-scores and analytical methods) selected from the report of the 5thILC on PAHs in edible oil [32]

515 used for checking the applicability of SRDs.

	L1		L2		L3		L4		L5		L6		L7		L8		L9		L10		L11		L12		L13		As.
	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg	z	µg/kg
5MC	1.20	0.35	1.61	2.02	1.11	0	1.01	-0.43	1.28	0.67	1.72	2.47	0.87	-1	1.60	1.98	1.50	1.57	1.30	0.75	1.13	0.06	1.30	0.75	0.98	-0.55	1.1
BAA	2.40	0.04	2.75	0.7	2.05	-0.63	2.48	0.19	2.81	0.82	1.57	-1.55	2.13	-0.48	2.60	0.42	3.80	2.71	2.30	-0.15	2.72	0.65	2.70	0.61	2.13	-0.48	2.4
BAP	2.90	-0.13	3.43	0.73	3.28	0.48	3.06	0.13	4.22	2.01	2.61	-0.61	2.40	-0.95	3.30	0.52	2.90	-0.13	2.80	-0.3	3.13	0.24	3.20	0.35	2.72	-0.43	3.0
BBF	5.20	-0.2	6.04	0.5	5.44	0	5.52	0.06	4.95	-0.41	4.64	-0.67	4.20	-1.04	7.00	1.3	6.20	0.63	5.40	-0.04	5.66	0.18	4.70	-0.62	5.16	-0.24	5.4
BCL	2.20	1.01	1.90	0.25	1.89	0.23	1.91	0.28	2.89	2.75	1.99	0.48	1.40	-1.01	2.00	0.5	1.80	0	1.30	-1.26	1.83	0.07	2.10	0.76	1.60	-0.51	1.8
BGP	6.10	-0.04	6.44	0.21	6.71	0.41	6.44	0.21	8.96	2.07	6.52	0.27	5.80	-0.26	6.80	0.48	5.50	-0.48	5.80	-0.26	6.58	0.31	6.30	0.11	5.97	-0.14	6.2
BJF	1.40	-0.07	1.78	1.15	1.78	1.15	1.49	0.22	11.73	32.99	1.40	-0.1	1.47	0.16	0.50	-2.95	1.90	1.53	1.70	0.89	1.01	-1.31	1.80	1.21	1.53	0.35	1.4
BKF	8.20	-0.03	10.08	1.01	8.65	0.22	8.23	-0.01	1.66	-3.63	8.53	0.16	6.00	-1.24	9.90	0.91	8.30	0.03	8.80	0.31	8.61	0.2	9.20	0.53	7.96	-0.16	8.2
CHR	3.70	0.45	4.12	1.02	3.60	0.31	3.58	0.29	4.21	1.14	4.17	1.08	3.33	-0.05	4.30	1.26	6.60	4.36	3.30	-0.09	3.87	0.68	4.40	1.39	3.28	-0.12	3.4
CPP	8.60	0.55	8.98	0.78	8.07	0.24	8.28	0.36	2.84	-2.86	11.00	1.98	5.33	-1.39	8.70	0.61	13.20	3.28	6.20	-0.87	7.17	-0.29	6.80	-0.51	6.20	-0.87	7.7
DEP	0.80	-0.97	1.02	0.02	1.03	0.06	0.78	-1.06	1.83	3.64	0.49	-2.35	1.00	-0.07	0.90	-0.52	0.80	-0.97	0.80	-0.97	0.78	-1.06	1.20	0.82	0.85	-0.75	1.0
DHA	4.90	1.33	5.17	1.65	5.08	1.54	4.83	1.24	4.98	1.42	8.00	5.05	4.00	0.25	5.70	2.29	4.00	0.25	4.40	0.73	5.05	1.51	4.80	1.21	4.52	0.87	3.8
DHP	2.10	-0.67	2.83	0.66	2.78	0.57	2.23	-0.44	4.60	3.92	2.20	-0.5	1.87	-1.1	3.00	0.98	2.80	0.61	2.10	-0.68	1.95	-0.96	4.20	3.18	2.11	-0.66	2.5
DIP	9.10	-0.3	10.69	0.44	10.60	0.4	9.31	-0.21	33.59	11.11	11.92	1.01	6.67	-1.44	10.30	0.26	11.20	0.68	10.30	0.26	9.41	-0.16	11.80	0.96	9.30	-0.21	9.8
DLP	1.60	0.43	1.77	0.95	1.51	0.15	1.13	-1.03	15.81	44.57	1.80	1.05	1.20	-0.82	1.70	0.74	1.70	0.74	1.60	0.43	1.41	-0.16	1.60	0.43	1.37	-0.29	1.5
ICP	3.40	-0.45	4.27	0.6	4.34	0.68	3.82	0.05	4.61	1.01	1.97	-2.17	3.53	-0.29	3.50	-0.33	3.40	-0.45	3.80	0.03	3.81	0.04	4.10	0.39	3.35	-0.51	3.8
% of acceptable results (z ≤2)	100		94		100		100		31		75		100		88		81		100		100		94		100		
Method used for PAHs analysis	SAP + LLE + GC-MS		SEC + LC-FLU		SEC + LC-FLU		SPE + GC-MS/MS		LLE + SPE + GC-MS		SAP + LLE + LC-FLU		LLE + SEC + GC-MS		SEC + LC-MS		SAP + SPE + GC-MS		SAP + SPE + GC-MS		SAP + LLE + GC-MS		SPE + SPE + GC-MS/MS		LLE + SEC + GC-MS		

516 SAP: saponification; LLE: liquid-liquid extraction; SEC: size exclusion chromatography; SPE: solid phase extraction; GC-MS: gas chromatography with mass spectrometry; LC-FLU: liquid chromatography with fluorescence detection; GC-MS/MS: gas chromatography with tandem mass spectrometry; LC-MS: liquid chromatography with mass spectrometry

518 Figure captions

519 Figure 1

520 PCA score plots PC1 vs. PC2 for the sets consisted of (a) the reported and assigned values

521 (“OIL+As”), and (b) the absolute values of z -scores (“Z-SCORE”) for the laboratories (methods)

522 (L1, L2, ..., L13)

523

524 Figure 2

525 The dendrogram of the laboratories according to a) the reported and assigned values (“OIL+As”

526 data set), and b) the absolute values of z -scores (“Z-SCORE” set)

527

528 Figure 3

529 Box and whisker plots of the absolute z -scores calculated for the laboratories (methods) (L1, L2, ...,

530 L13) a) all 13 laboratories included in the “Z-SCORE” formed in this study, b) after excluding L5

531 as an outlier.

532

533 Figure 4

534 Line plots for SRD rankings: “OIL” set, reference: assigned value (full circles, blue); “OIL+As” set,

535 reference: averages (full boxes, red); “OIL+As-OUT” set, reference: averages (full rhombuses,

536 green); “OIL+As”, reference: medians (full triangles, pink); normalized sum of squared z -scores,

537 SZ2norm (black full circles, dotted line)

538

539 Figure 5

540 SRD ranking with CRNN validation of 13 laboratories for a) “OIL” set, b) “OIL+As” set. The Y

541 left-hand side-axis and X-axis are SRD values scaled between 0 and 100. The Y right-hand side-

542 axis represents relative frequencies of the theoretical distribution for ranking random numbers.

543 Statistical characteristics of this distribution (CRRN procedure) are defined by the first icosaille
544 (5%), XX1; the first quartile, Q1; median, Med; the last quartile, Q3; the last icosaille (95%), XX19.

545

546 Figure 6

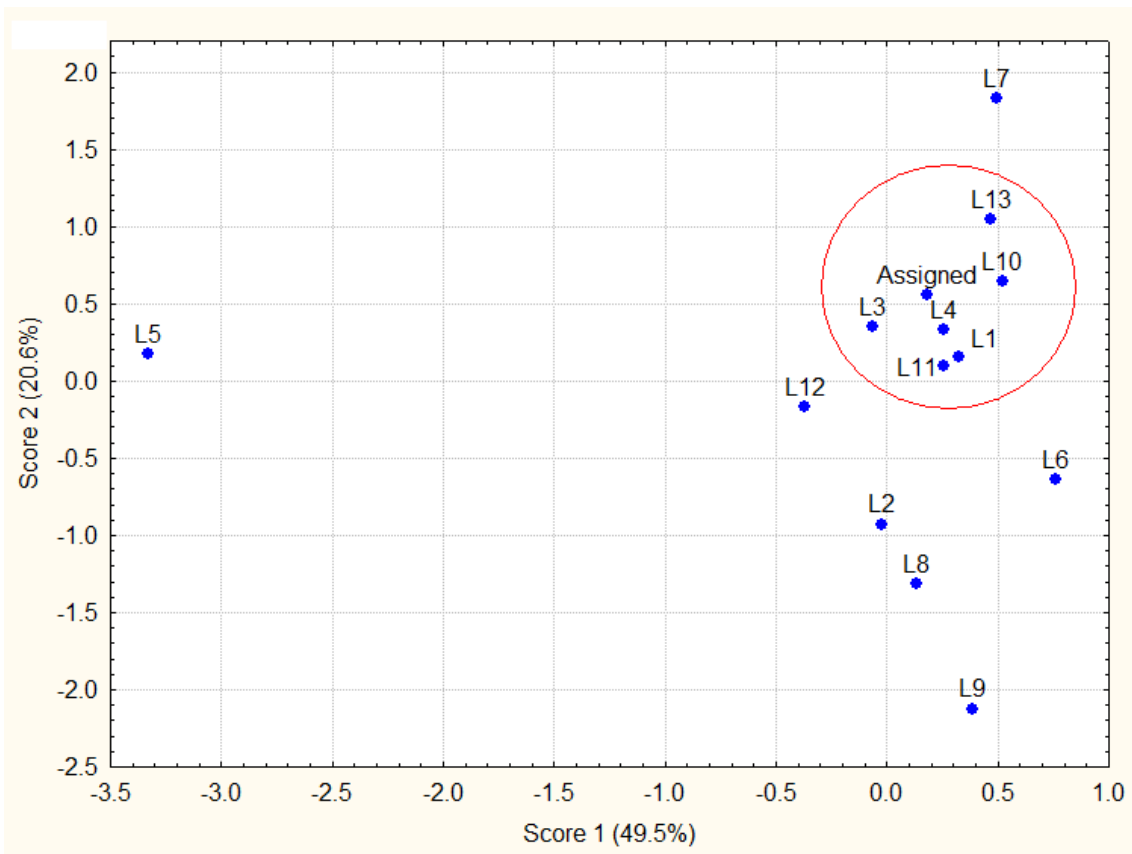
547 SRD ranking with CRNN validation of 13 laboratories according to the absolute values of z -scores
548 (“Z-SCORE” set) calculated according the contents of EU 15+1 PAHs reported during the 5thILC
549 on PAHs organized by IRMM, Geel, Belgium [32]. The Y left-hand side-axis and X-axis are SRD
550 values scaled between 0 and 100. The Y right-hand side-axis represents relative frequencies of the
551 theoretical distribution for ranking random numbers. Statistical characteristics of this distribution
552 (CRRN procedure) are defined by the first icosaille (5%), XX1; the first quartile, Q1; median, Med;
553 the last quartile, Q3; the last icosaille (95%), XX19.

554

555 Figure 7

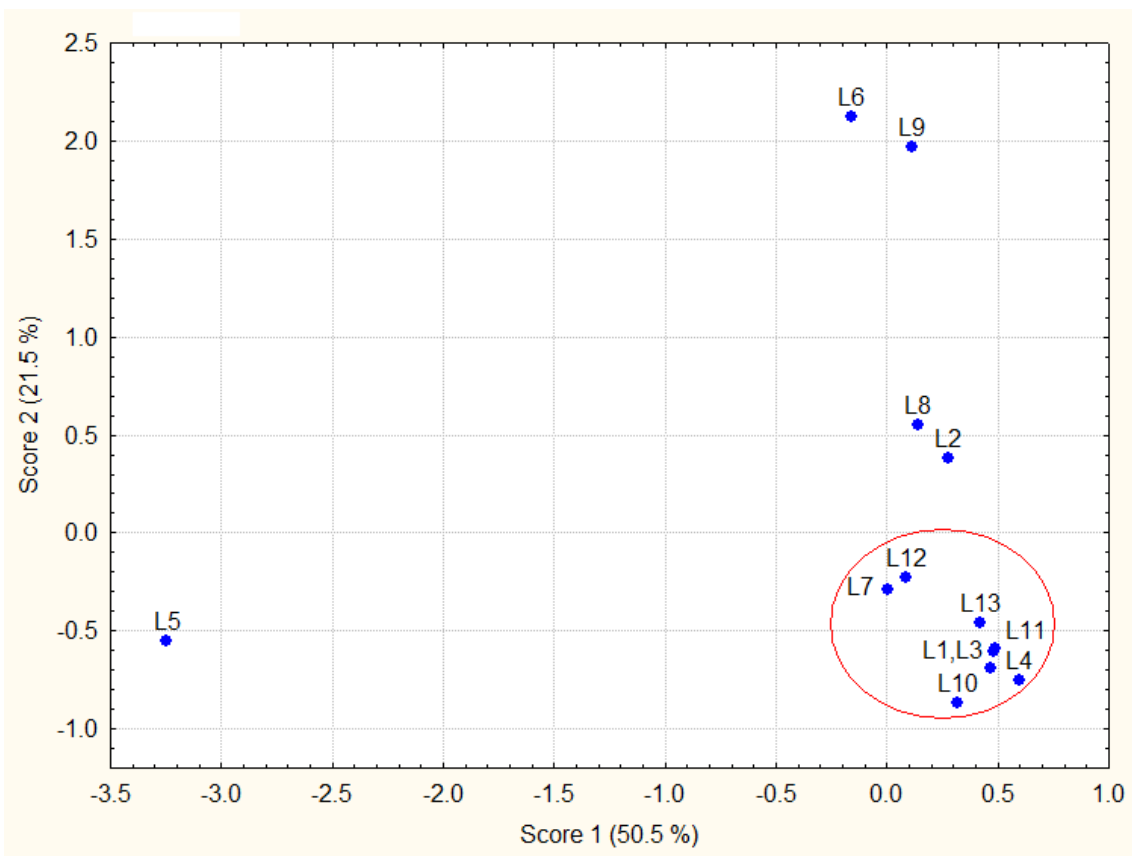
556 Box and whisker plot of the original PAH concentrations (a); box and whisker plot of sum of
557 ranking difference values obtained from a seven segments cross-validation (b).

558



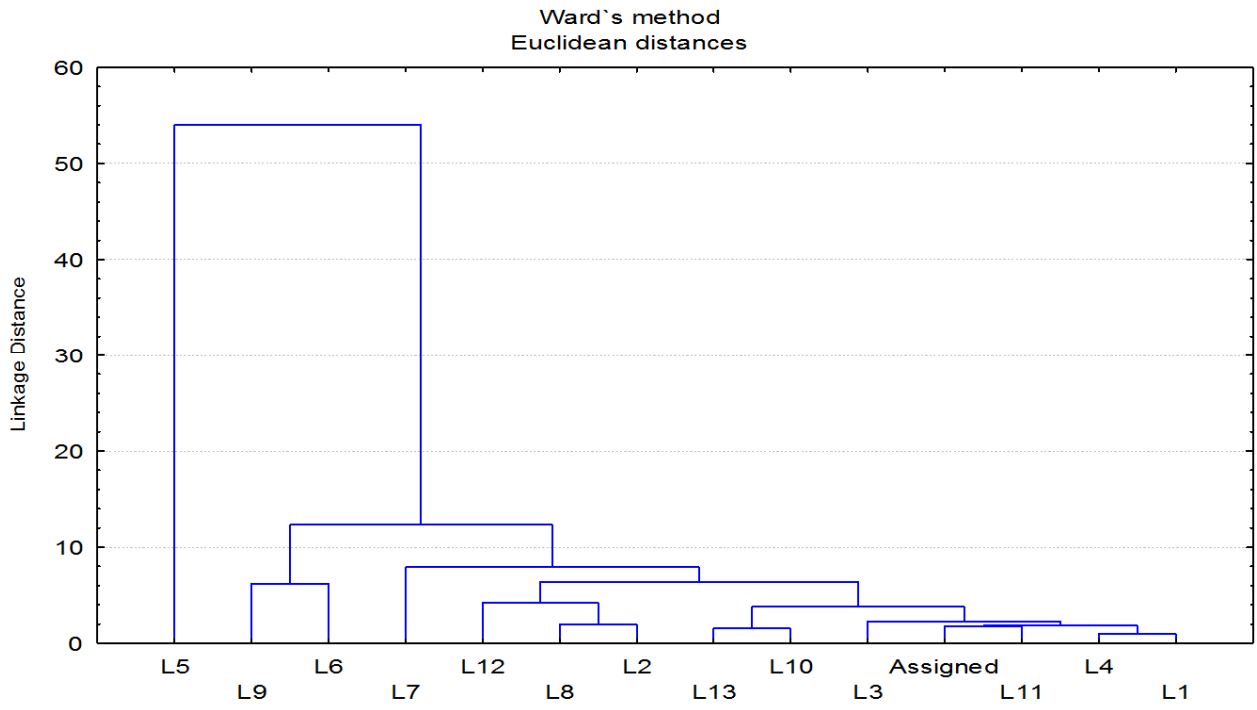
559

560 Figure 1a



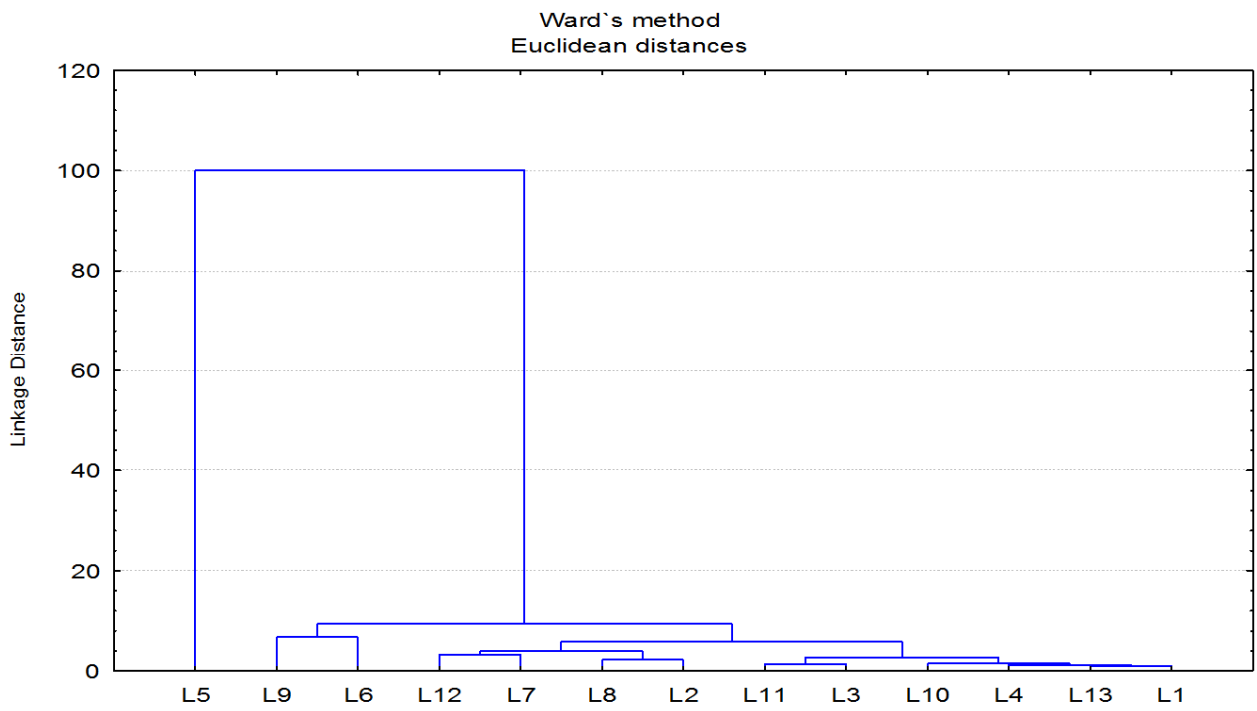
561

562 Figure 1b



563

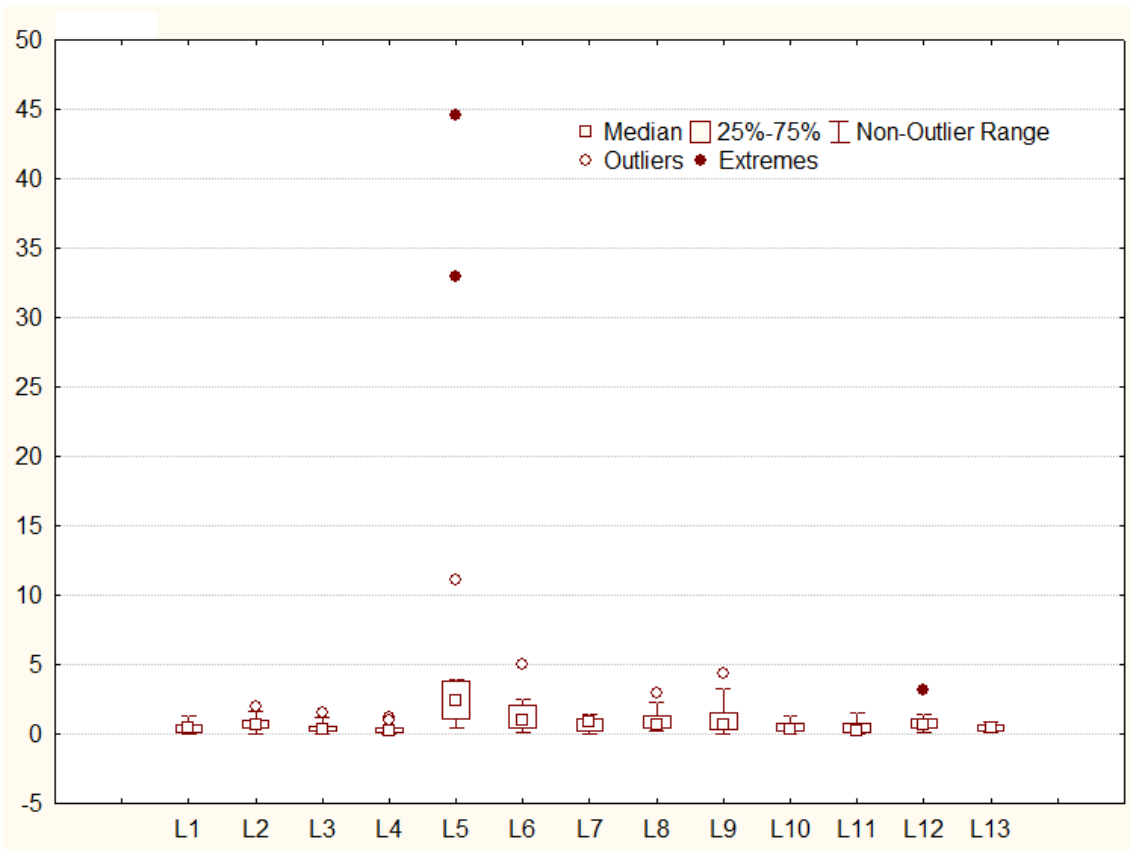
564 Figure 2a



565

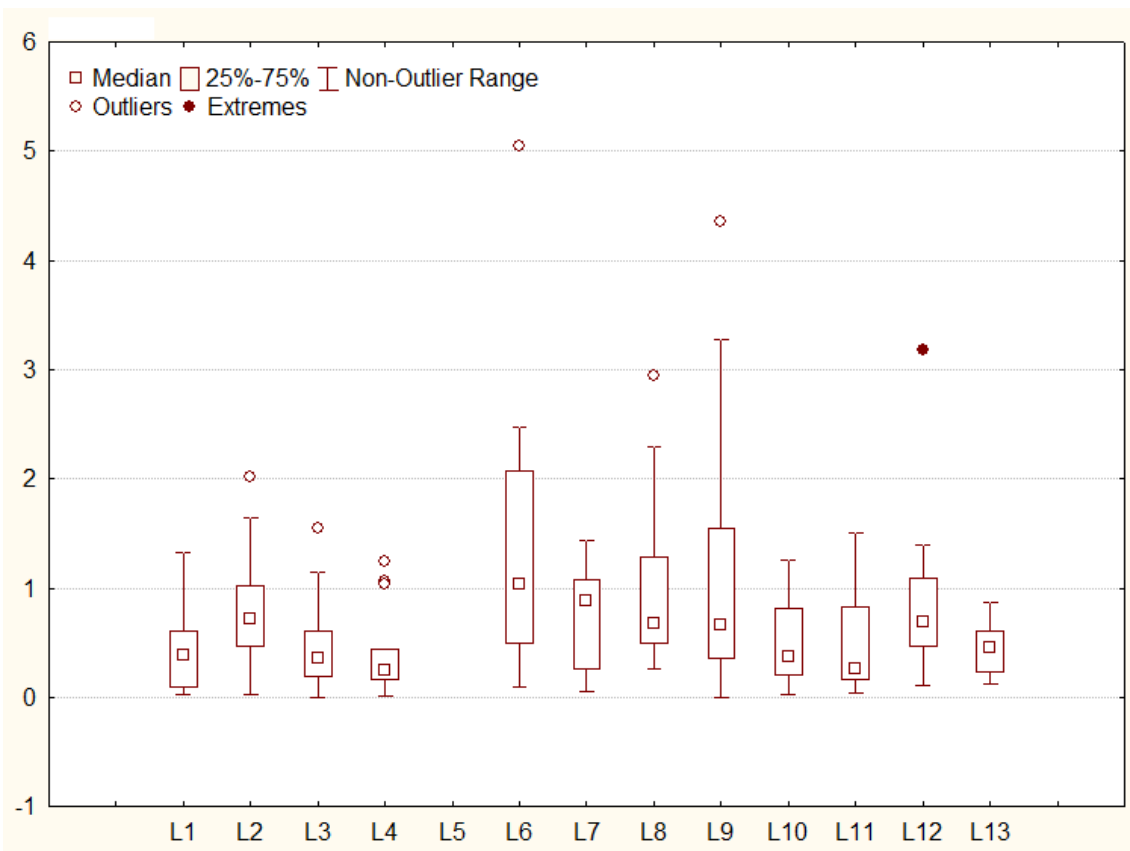
566 Figure 2b

567



568

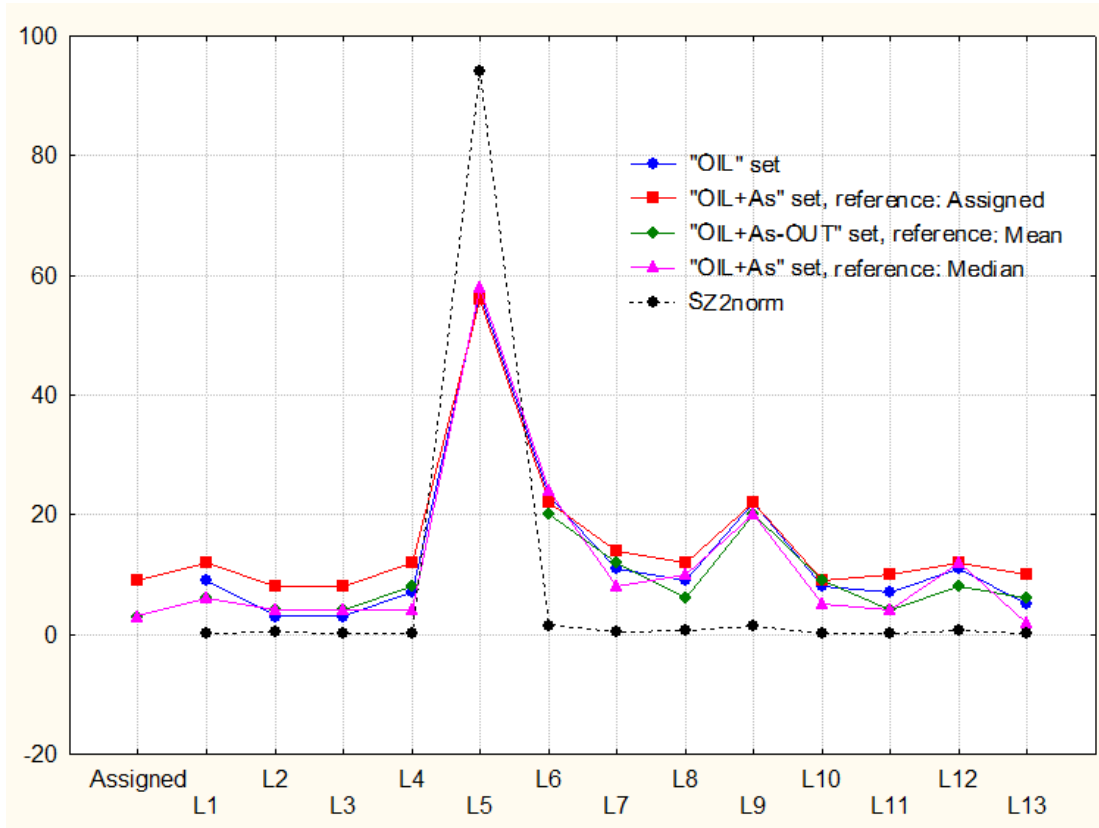
569 Figure 3a



570

571 Figure 3b

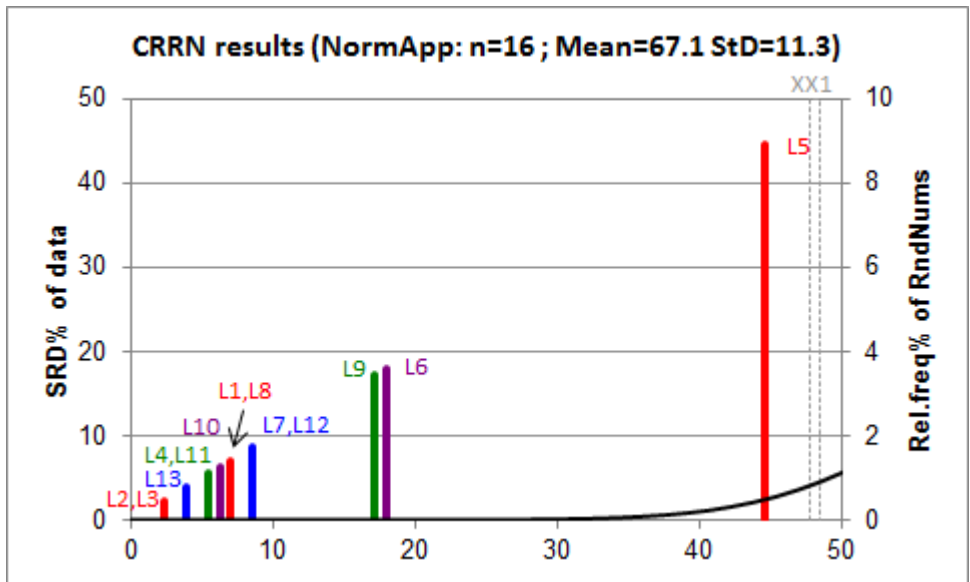
572



573

574 Figure 4

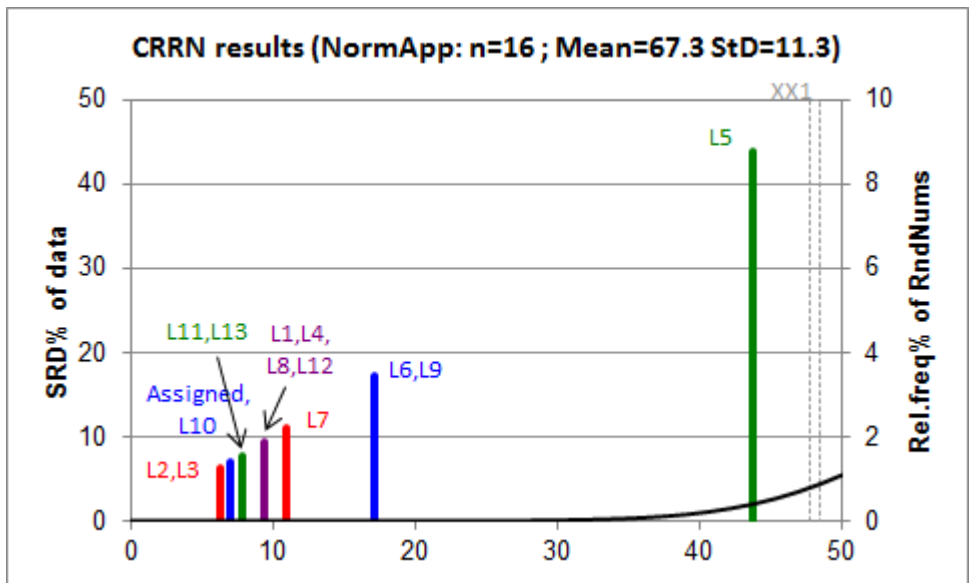
575



576

577 Figure 5a

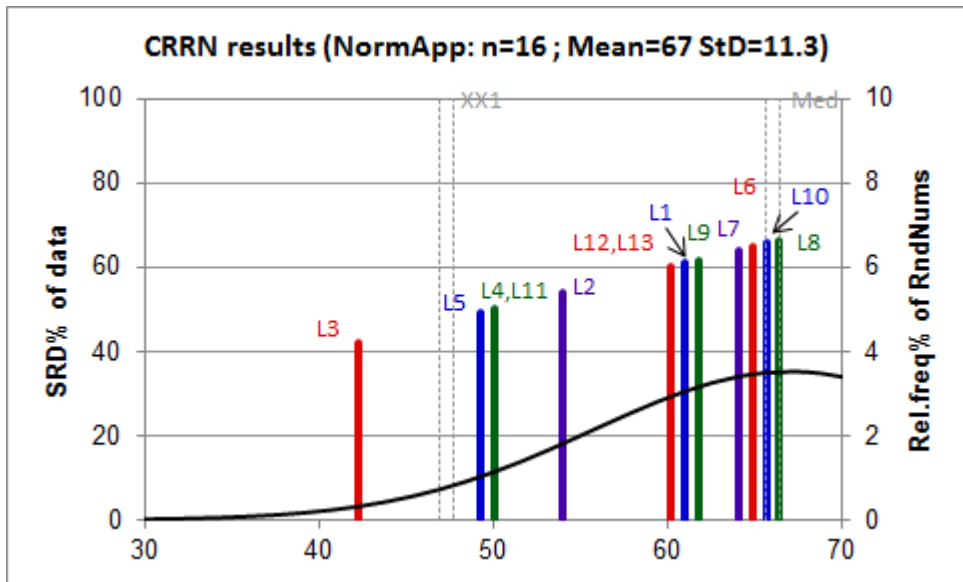
578



579

580 Figure 5b

581



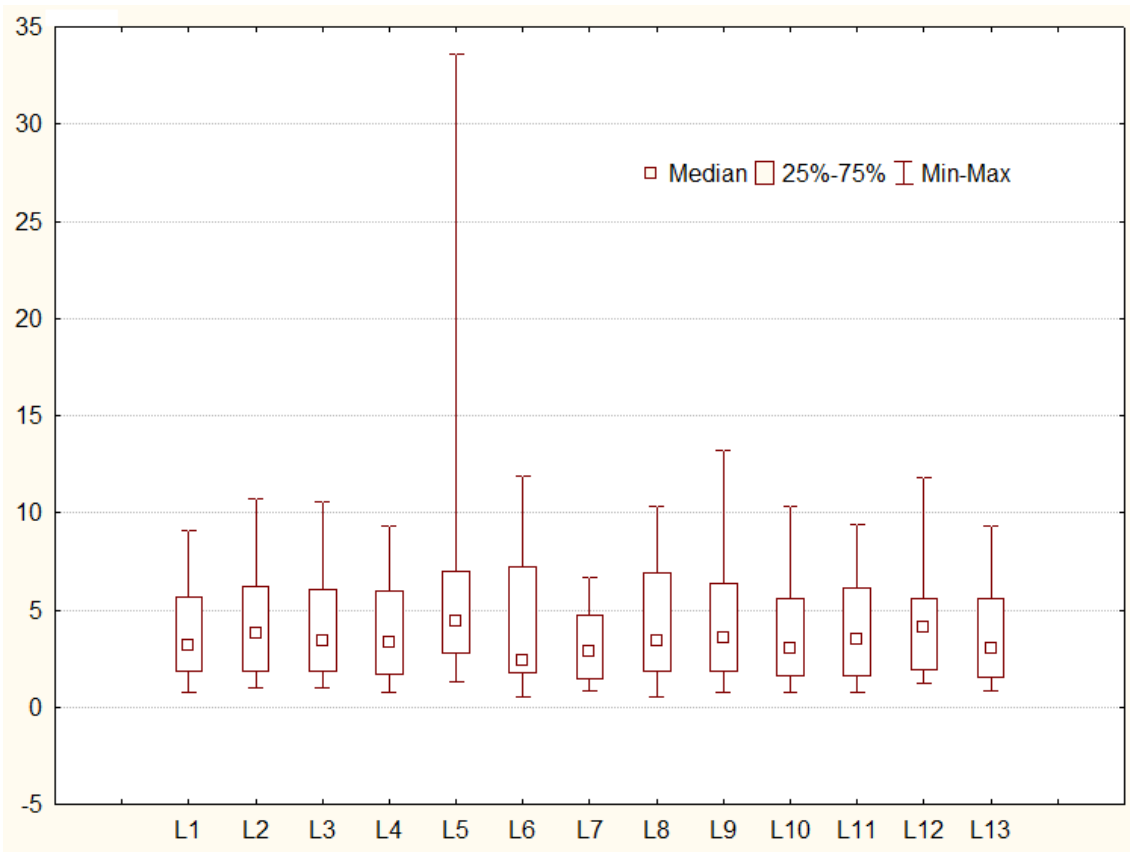
582

583 Figure 6

584

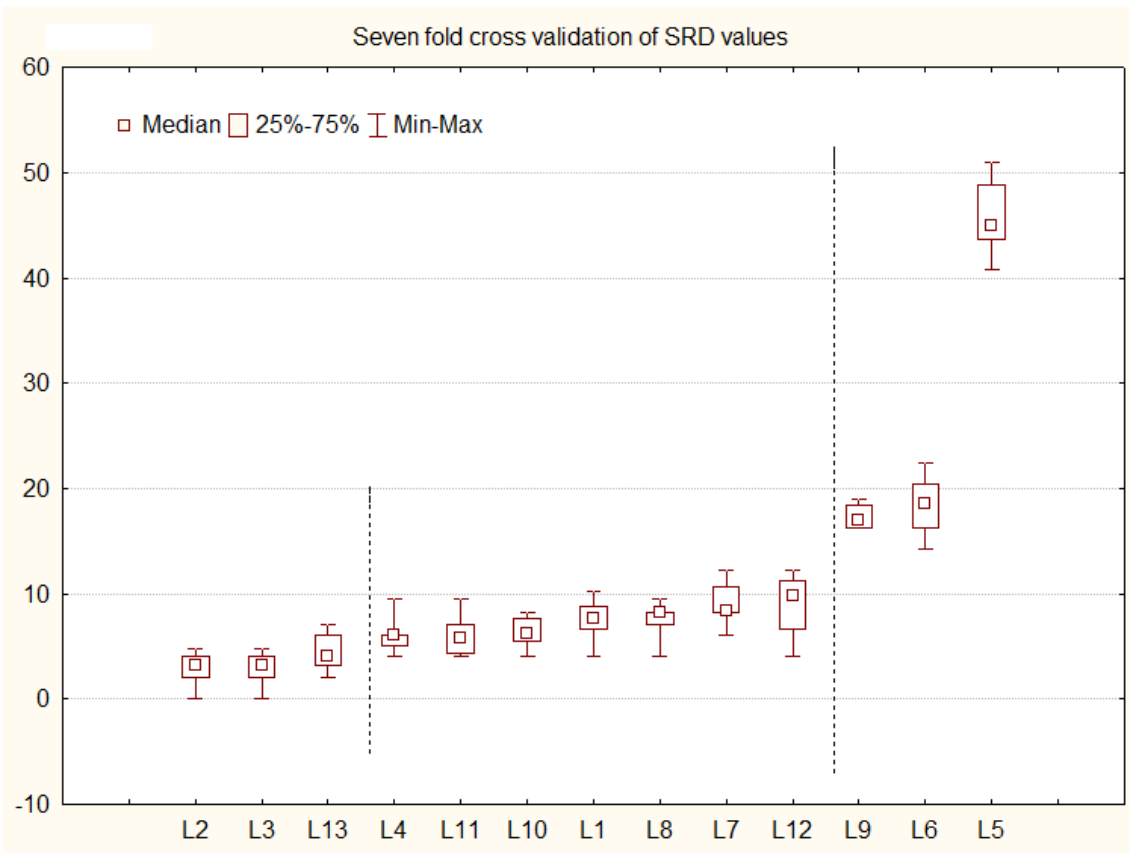
585

586



587

588 Figure 7a



589

590 Figure 7b