# EVALUATION OF DISTANCE-BASED APPROACHES FOR FORENSIC COMPARISON: APPLICATION TO HAND ODOR EVIDENCE

**SCHOLARONE™**
Manuscripts

1
2
3    **EVALUATION OF DISTANCE-BASED APPROACHES FOR FORENSIC**

4    **COMPARISON: APPLICATION TO HAND ODOR EVIDENCE**

3    **ABSTRACT**

4    The issue of distinguishing between the same-source and different-source hypotheses based

5    on various types of traces is a generic problem in forensic science. This problem is often

6    tackled with Bayesian approaches, which are able to provide a likelihood ratio that quantifies

7    the relative strengths of evidence supporting the two competing hypotheses. Here, we focus

8    on distance-based approaches, whose robustness and especially capacity to deal with high-

9    dimensional evidence are very different, and need to be evaluated and optimized.

10   A unified framework for direct methods based on estimating the likelihoods of the distance

11   between traces under the two competing hypotheses, and indirect methods using logistic

12   regression to discriminate between same-source and different-source distance distributions,

13   is presented. They are compared in terms of sensitivity, specificity and robustness, with and

14   without dimensionality reduction, with and without feature selection, on the example of hand

15   odor evidence. Empirical evaluations on a large panel of 534 subjects show the significant

16   superiority of the indirect methods, especially without dimensionality reduction.

17   **KEYWORDS**

18   Bayesian inference; dissimilarity measure; forensic science; human hand odor; likelihood

19   ratio; logistic regression.

20   **HIGHLIGHTS**

21   • Direct and indirect distance-based likelihood ratio estimation methods for forensic

22   comparison are investigated

23   • They are applied to high-dimensional evidence consisting of hand odor traces

24   • Their robustness, AUC, sensitivity and specificity are evaluated on a panel of 534 subjects

25   • Indirect methods based on logistic regression outperform direct ones and are more robust

26   • Indirect methods using a vectorial distance outperform those using a scalar one, with and

27   without feature selection

28

29

## 1. INTRODUCTION

A generic problem in forensic science is to decide whether a trace of an unknown source, often drawn from a crime scene, and a trace from a known source, stem from the same source, a person or a firearm for example. If the source is a person, the traces might be biometric such as a DNA profile [Aitken and Taroni 2004] [Puch-Solis et al. 2012], fingerprints [Neuman et al. 2012], a voice [Morrison 2011], an olfactory profile [Cuzuel et al. 2017], or they might consist of footwear impressions [Tang and Srihari 2014bis], handwriting [Tang and Srihari 2014], etc. If the source is a firearm, the traces may be features such as striations and impressions of a bullet or of a cartridge case [Mattijssen et al. 2020].

The most common approach to this problem is to estimate a likelihood ratio (LR), i.e. the ratio of the joint probability of occurrence of the two traces under the hypothesis that they arose from the same source and under the hypothesis that they arose from different sources. A convenient solution is to replace the joint probability of the traces by the probability of a distance between the two traces quantifying their dissimilarity [Neuman et al. 2007, Riva and Champod 2014, Tang and Srihari 2014bis, Ali et al. 2015, Muehlethaler et al. 2016, Riva et al. 2020, Vergeer et al. 2020]. If, as is most often the case, the distance is scalar, there is an important loss of information. Thus, we choose to focus on distance-based methods, but with the possibility to use a vectorial distance between traces.

Furthermore, the distance-based LR estimate can be obtained either directly, by estimating the distance likelihoods under the two hypotheses, or indirectly, by first using logistic regression to discriminate between same-source and different source distance distributions, and then Bayes' formula to infer the LR. Whilst the direct method is more flexible, the indirect method is more robust and quite natural in machine learning [Bishop 2006, Hastie et al. 2009]. It is sometimes advocated for in the forensic context, for the same reasons and also because it enables score calibration and fusion with minimal mathematical complexity [Enzinger et al. 2016, Morrison 2013]. Here, we show that the indirect method also enables the use of a vectorial distance, thus preventing the severe information loss suffered by scalar distance approaches. We discuss the direct and indirect methods in terms of robustness and ability to handle high dimensional evidence, with or without dimensionality reduction, and with or

59  without feature selection. We evaluate them in terms of sensitivity, specificity and robustness

60  on the example of traces consisting in a hand odor profile.

61  **2. MATERIALS AND METHODS**

62  **2.1 Problem statement**

63  The aim is, given the evidence consisting of a pair of traces (e.g. two olfactory profiles), to

64  decide whether these traces have the same source (e.g. the same person) or not. In the

65  following, $H_{ss}$ refers to the hypothesis that the two traces stem from the same source, and $H_{ds}$

66  to the alternative hypothesis that they stem from different sources. Given the *a priori*

67  probabilities $P(H_{ss})$ and $P(H_{ds})$, the Bayesian formula yields the posterior probability of $H_{ss}$

68  given the evidence E:

$$P(H_{ss} \mid E) = \frac{f(E \mid H_{ss})P(H_{ss})}{f(E \mid H_{ss})P(H_{ss}) + f(E \mid H_{ds})P(H_{ds})}$$

69                                                                                                          (1)

70  where $f(E \mid H_{ss})$ and $f(E \mid H_{ds})$ are the distributions of the evidence under $H_{ss}$ and $H_{ds}$, or

71  likelihoods. Jeffreys developed an absolute scale to evaluate the degree of confidence in the

72  same-source hypothesis outside a decisional framework based on the posterior probability of

73  $H_{ss}$ using the LR [Jeffreys 1939, Robert 2001] defined as:

$$LR(E) = \frac{f(E \mid H_{ss})}{f(E \mid H_{ds})}$$

74                                                                                                          (2)

75  which is independent from the *a priori* probabilities. In fact, the observed evidence E consists

76  of the two traces which are represented by n-dimensional vectors (whose components are

77  the amounts of each odor compound). Since we focus on distance-based methods, the

78  information contained in these two vectors is transformed into a distance or dissimilarity

79  measure, which can be either a scalar or a n-vector (a distance for each feature of the trace,

80  here for each odor compound). In the following, this distance between the two traces will be

81  denoted by *d*.

82  **2.2 Candidate methods**

83  The LR can be obtained either directly, i.e. by estimating the likelihoods of the distance under

84  the two competing hypotheses, or indirectly, i.e. using first logistic regression to discriminate

85  between same-source and different source distance distributions, and then formulas (1) and

86  (2) to infer the LR.

### 2.2.1 Direct methods

87

88  For direct methods, we need to be able to evaluate $f(d|H_{ss})$ and $f(d|H_{ds})$ whatever the value

89  of d. For this purpose, part of the available dataset can be used to build pairs of traces of the

90  two types, same-source and different-source pairs. If d is scalar, or of dimension 2 or 3 at

91  most, the two empirical distributions can be fitted, with Gaussian mixtures, for example,

92  leading to parametric estimates of $f(d|H_{ss})$ and $f(d|H_{ds})$ [Mattijssen et al. 2020, Riva and

93  Champod 2014, Riva et al. 2020]. In the case of many features, a fit of each component of d

94  can be performed in the same way, and the overall likelihoods can be approximated through

95  the product of the likelihoods in each dimension, leading to the naïve Bayes classifier. To be

96  successful, the latter approach necessitates however that the features are not overly

97  correlated.

### 2.2.2. Indirect methods

98

99  The aim of these method is to build a discriminative model of the boundary between the two

100  categories of pairs (same-source and different-source pairs) rather than a generative model

101  explicitly parameterizing the distributions in the two categories. It is well known that, under

102  the hypothesis of single-Gaussian distributions with the same variance under $H_{ss}$ and $H_{ds}$ in

103  the scalar case, or same covariance matrix in the multidimensional case, the posterior

104  probability of $H_{ss}$ takes the form of a sigmoidal curve [Bishop 2006, Hastie et al. 2009], hence

105  the motivation for a logistic regression approach. Despite its result being a discriminative

106  model, it enables to calculate the posterior probability of Equation (1) as well as the LR of

107  Equation (2). As a matter of fact, the logistic regression model with parameters $\theta = [a^T\ b]^T$ has

108  output:

$$r(d,\theta) = \frac{1}{1 + \exp\left(-(a^T d + b)\right)}$$

109                                                                                                  (3)

110  where $b$ is a scalar, and $a$ is either a scalar in the case of a scalar distance, or otherwise a n-

111  vector (of the dimension of the evidence). If the proportions of the same-source and different-

112  source categories in the calibration set are denoted by $f_{ss}$ and $f_{ds}$, r(d, θ) approximates:

$$\frac{f(d|H_{ss})f_{ss}}{f(d|H_{ss})f_{ss} + f(d|H_{ds})f_{ds}}$$

113                                                                                                  (4)

114  Thus, the posterior probability for *a priori* probabilities $P(H_{ss})$ and $P(H_{ds})$ can be retrieved with:

$$P(H_{ss} \mid d) = \cfrac{1}{1 + \exp\left(-(a^T d + b)\right) \cfrac{f_{ss}}{f_{ds}} \cfrac{P(H_{ds})}{P(H_{ss})}}$$

(5)

116 and the likelihood ratio with:

$$LR(d) = \exp\left(a^T d + b\right) \frac{f_{ds}}{f_{ss}}$$

(6)

118 **2.3. Pros and cons**

119  The indirect method offers several advantages:

120 - it spares the necessity to fit the likelihoods,

121 - in the multi-dimensional case, contrary to the naïve Bayes classifier, the independence

122 assumption is not necessary, because the logistic regression automatically takes care of the

123 correlation between features,

124 - by construction, the log LR is defined by a hyperplane, and thus robust with respect to the

125 equal variance assumption, and to outliers or sparse data far from the boundary,

126 - in the forensic context, since the log LR is directly proportional to $a^T d + b$, the logistic allows

127 a convenient and interpretable calibration of the dissimilarity score d, and a fusion of scores

128 in the multidimensional case [Morrison 2013].

129 On the other hand, the indirect method might suffer from:

130 - a reduced flexibility since it amounts to assume single-Gaussian distributions,

131 - a possibly important computation time in the case of high-dimensional evidence and of a

132 distance of the same dimension.

133 These advantages and disadvantages will be examined and discussed on the example of hand

134 odor evidence using a large panel of subjects.

135 **2.4. Dataset description**

136 A panel of 534 volunteers was set up which gathers 218 men and 316 women aged 7 to 94

137 years (median 28, interquartile interval [22 ; 48]), see Table 1 for the detailed composition in

138 terms of sex and age. Note that this composition does not aim at reflecting that of a precise

139 target population, such as one which is more likely to commit a crime, but to be as

140 representative of the diversity of odors as possible. As a matter of fact, criminal investigations

141 also often necessitate to look for victims, or to discriminate between traces from different

142 people present at a crime scene, including those of victims or witnesses, who might be women

143 as well as men, children or seniors as well as middle-aged adults. All data were completely

144 anonymized prior to analysis, and no personal information was stored.

145 The goal was here to identify the subjects by their hand odor, whose volatile profile was shown

146 to display a between-subject variability which is sufficient for differentiation [Curran et al.

147 2010, Cuzuel et al. 2017]. Also, in the forensic context, the hands have the advantage to be

148 more likely to be directly in contact with objects at a crime scene, and to be easier to sample

149 during a police interrogation.

150 The volatile profiles were obtained by a direct sampling procedure using identical sample

151 collection kits of 4 small polymer strands that the subjects were asked to rub together in their

152 hands for 15 minutes. The polymer strands were thermodesorbed, and the concentrated

153 substances were separated by comprehensive bidimensional gas chromatography (GCxGC)

154 coupled with mass spectrometry (MS). The sampling method and the optimization of the

155 GCxGC-MS analysis were extensively described in [Cuzuel et al. 2017bis, Cuzuel et al. 2018].

156 Data were acquired, converted to .mzXML files with GC Real Time Analysis 4.20 (Shimadzu

157 software), and then processed with MatlabTM (Natick, MA, USA) version 9.6.0.1150989

158 (R2019a), its Statistics and Machine Learning Toolbox version 11.5 and its Bioinformatics

159 Toolbox version 4.12.

160 Using a "home-made" Matlab script [Cuzuel 2017], the preliminary manual processing of 25

161 chromatograms obtained on 3 subjects between 23 and 26 years old of both genders sampled

162 several times at different time instants enabled us to draw up a first list of several hundreds

163 of peaks. A library was built to store their retention times, their linear retention index, their

164 mass spectrum, and the name of the corresponding compound when it could be identified

165 using the NIST library. Indeed, if the availability of its mass spectrum is compulsory, a

166 compound does not need to be formally identified for the comparison of chromatograms. We

167 also checked whether compounds described in the literature as constituents of the human

168 hand odor were present in this library, otherwise they were included. The library was then

169 continuously enriched as the panel was increased with compounds potentially relevant to

170 human hand odor because of their empirical frequency in new samples. This work led us to a

171 customized library of 741 compounds, which were looked for in each chromatogram. As a

172 result, each sample was characterized by the peak area of 741 compounds.

173 In order to compensate for uncontrolled variations of the total area of the chromatograms,

174 the sum of these areas was normalized in logarithmic scale to unit value, see Table 2 for a

175  comparison of the reproducibility of the data without normalization, and normalization in

176  scalar and logarithmic scales. Not knowing whether all 741 compounds are really relevant for

177  identification (the median frequency of presence of the compounds across all the samples is

178  of 40.2 %, quartiles [12.8 %; 77.6 %]), such a normalization might be questionable. Thus, the

179  possibility to avoid the problem by working on the binarized areas, i.e. 1 if the compound is

180  present, or 0 if it is absent from the sample, was also investigated. Also, this approach might

181  be of interest for forensic identification problems dealing with intrinsically binary features,

182  such as gradient, structural and concavity (GSC) binary features in handwriting identification

183  [Srihari et al. 2008]. In the following, we refer to these traces of 741 features, continuous or

184  binarized, as "odor traces".

185  As stated above, the subjects were sampled in quadruplicate, but due to unavoidable mishaps

186  with some samples (like accidently dropping a polymer on the floor during sampling) and to

187  chromatographic problems (such as failures of the cryogenic modulator), 1690 odor traces

188  were obtained for the 534 subjects (44 were sampled once, 77 twice, 160 three times, and

189  the remaining 253 subjects four times, leading to an average of 3.2 odor traces per subject).

190  This data set was split into a calibration set for training and validation, and an independent

191  test set for performance estimation. The split was made so as to respect the gender

192  proportions, with subject of all ages in the two sets, and odor traces of the same subject being

193  put in the same set. As a result, the calibration set comprises 412 subjects and their 1 299

194  odor traces (corresponding to 1 594 $H_{ss}$ and 841 457 $H_{ds}$ pairs), and the test set comprises the

195  remaining 122 subjects and their 391 odor traces (leading to 481 $H_{ss}$ and 75 764 $H_{ds}$ pairs). The

196  way the odor traces distribute between calibration and test set can be grasped through the

197  Principal Component Analysis (PCA) of Figure 1.

198  **2.5. Implementation**

199  Three methods are implemented:

200  1) the direct method using a scalar distance between odor traces,

201  2) the indirect method using a scalar distance between odor traces,

202  3) the indirect method using a vectorial distance, i.e. a distance on each odor compound.

203  Note that, given the large dimension of the problem (n=741) and the known correlations

204  between features, we did not attempt to implement the direct method using a vectorial

205  distance (i.e. a scalar distance on each feature and the naïve Bayes classifier).

### 2.5.1. Distances between two odor traces

Concerning the choice of the scalar distance, correlation-based distances are robust with respect to shifts and linear transformations of the features, and since Spearman's correlation coefficient is able to capture a monotonic nonlinear association, as opposed to Pearson's linear correlation coefficient [Daniels 1944], the Spearman correlation based distance is also expected to be more robust with respect to nonlinear variations of peak areas. Since this proved to be true in a previous study where Euclidian distance, Pearson correlation and Spearman correlation based distances were compared [Cuzuel et al. 2018bis], we restrict here to Spearman's correlation based distance for the direct and indirect methods using a scalar distance. The chosen vectorial distance for the third method is simply the vector of the absolute differences between feature values.

### 2.5.2. Estimation of the likelihoods for the direct method

The calibration set was used to build pairs of odor traces of same and different sources, and to compute their distances. The empirical densities were fitted with a two-Gaussian mixture distribution, using Matlab's function "fitgmdist", leading to estimates of the likelihoods $f(d|H_{ss})$ and $f(d|H_{ds})$.

### 2.5.3. Estimation of the logistic model for the indirect methods

The logistic regression model of Equation (3) with parameters $\theta = [a^T \, b]^T$ was fitted to minimize the cross-entropy cost function using Matlab's function "glmfit".

### 2.5.4. Likelihood ratio and performance estimation

For the direct method, the LR was evaluated using the estimates of the likelihoods $f(d|H_{ss})$ and $f(d|H_{ds})$ and Equation (2), as a function of the distance d. For the indirect methods, the LR was obtained from the fitted logistic regression and Equation (6), and plotted as a function of the distance d or of the score $a^T d + b$, depending on d being scalar of vectorial.

Since there is no true reference for the LR, the performance of the different methods was evaluated by estimating the posterior probability $P(H_{ss}|d)$ according to Equations (1) and (5) for the direct and indirect methods respectively, and by performing a binary classification using equal prior probabilities ($P(H_{ss}) = P(H_{ds}) = 0.5$). Varying the decision threshold on $P(H_{ss}|d)$, the sensitivity and the specificity were estimated on the calibration and test sets, and used to compute the corresponding areas under the receiver operating characteristic

236 "ROC" curve (AUC) [Hanley and McNeil 1982]. The performance was further characterized by

237 the sensitivity and specificity maximizing Youden's index [Youden 1950], i.e. their sum.

### 2.5.5. Feature selection

239 In a previous study [Cuzuel et al. 2018bis], improved results were obtained using feature

240 selection. Given the large number of features (odor compounds) and the large size of the data

241 set, an economic and robust filter approach to this selection was chosen. The idea is to retain

242 the features that contribute the most to the difference between densities under $H_{ss}$ and $H_{ds}$.

243 For each feature, using the absolute values of the difference for the $H_{ss}$ and $H_{ds}$ pairs, we

244 computed Wilcoxon's non-parametric test statistic in the case of continuous features, and

245 Fisher's exact test statistic in the case of binarized features. Then, the features were ranked

246 in decreasing order of the one-sided p-value of the test (it is a one-sided test since smaller

247 differences between features under $H_{ss}$ than under $H_{ds}$ are sought for). The number of features

248 maximizing the AUC was estimated on the calibration set using 3-fold cross-validation. The

249 cross-validation partitions were randomly chosen with the constraint that the odor traces of

250 the same subject were put in the same partition. Note that cross-validation also enabled us to

251 estimate the uncertainty on the AUCs through the mean standard deviation on the three

252 partitions.

### 3. RESULTS AND DISCUSSION

254 The three methods are first evaluated using all the features of the odor traces (baseline

255 comparison) and then, the possibility to further improve their performance using feature

256 selection is investigated.

### 3.1. Baseline comparison of the three methods (without feature selection)

258 The results obtained with the three methods on the calibration and test sets are summarized

259 in Tables 3 and 4 for binarized and continuous features respectively. As a first remark, the

260 performance of the direct and indirect methods using a scalar distance in terms of AUC and of

261 sensitivity and specificity are almost identical, for both binarized and continuous features.

262 Thus, the higher flexibility of the direct method does not increase the performance. On the

263 contrary, its lack of robustness can be visualized on Figure 2 depicting the posterior probability

264 and the LR obtained with the binarized features: due to the larger variance of the likelihood

265 under H$_{ss}$, the posterior probability and the LR, instead of being monotonous, start to increase

266 with the distance at some point (d $\approx$ 0.7). Whereas whatever the situation with the indirect

267 method, posterior probability and LR always decrease with d, see Figure 3 depicting the

268 posterior probability and the LR obtained with indirect method, this time on the continuous

269 features.

270 Also noteworthy, the performance obtained with the indirect method using a vectorial

271 distance is significantly better than those of the methods working with a scalar distance: the

272 AUC on the calibration and test sets jumps from 91-92% to 97-98%, the standard deviation of

273 the AUC being estimated at 0.7% using 3-fold cross-validation on the calibration set. The

274 distributions of the score (a$^T$ d + b) resulting from the logistic regression, the regression itself,

275 the posterior probability and the LR are shown in Figure 4. The only drawback lies in the

276 increased, but perfectly tractable computational cost (10 minutes instead of a few seconds,

277 on a 4,2 GHz Intel Core i7).

278 Finally, the binarization of the features decreases the performance, but only marginally (the

279 AUC is decreased by $\approx$ 1%). Note that, in this precise case where the features quantify the

280 amount of odor compounds, this could be due to the fact that the normalization of the

281 compound proportion uses all these compounds whereas it is not known whether they are all

282 relevant. Note also that the normalization was improved by performing it in the logarithmic

283 scale rather than in the linear scale (the former improving the reproducibility, see Table 2),

284 with which continuous features did not outperform binarized features, as shown in a previous

285 study [Cuzuel et al. 2018bis]. Finally, other normalization methods specific to GCxCG-MS data

286 might advantageously be investigated [Chen et al. 2017], but are outside the scope of this

287 paper.

**3.2. Comparison of three methods with feature selection**

289 The number of selected features using the filter approach is reported in Tables 5 and 6,

290 together with the corresponding results on the calibration and test sets, for binarized and

291 continuous feature respectively.

292 Again, there is almost no difference in performance between the direct and indirect methods

293 with a scalar feature, be it on binarized or continuous features. In terms of AUC, the selection

294 is more efficient on binary features than on continuous ones (94.4% with selection instead of

295 91.5% without for binarized features, 93.5% instead of 93.0% for continuous features, on the

296  test set), with an important reduction of the number of the binarized features (267 instead of

297  741), and a moderate one for continuous features (535 instead of 741). Note also that in both

298  cases, this increased performance benefits the specificity, which is highly desirable in a

299  forensic application (it is crucial in this context not to reject the different-source hypothesis,

300  i.e. the defense hypothesis, when in fact it is true).

301  For both binarized and continuous features, the indirect method using a vectorial distance is

302  again significantly better than the previous ones (AUCs around 97-98% instead of 94-95% on

303  both calibration and test sets), with a similar number of selected features (440 for binarized

304  features, 500 for continuous ones). However, in both cases, the parsimony due to feature

305  selection does not increase the performance as compared to the baseline method, it is quasi-

306  identical with and without selection. In return, this testifies to a robustness of the indirect

307  method with respect to possibly irrelevant features. And of course, not to have to perform

308  the selection spares computation time.

### 3.3. Discussion of the choice of equal priors

310  In this manuscript, the methods are compared in terms of AUC, sensitivity and specificity. In

311  the case of the indirect methods, the regression being obtained by fitting a logistic function to

312  the data, whatever the prior probabilities $P(H_{ss})$ and $P(H_{ds})$, the posterior probability $P(H_{ss}|d)$

313  given by Equation (5) is also a logistic function. Thus, when the threshold on $P(H_{ss}|d)$ is varied

314  from 1 to 0, the same ROC curve is described, whose AUC only depends on the distance

315  distributions under $H_{ss}$ and $H_{ds}$: only the threshold maximizing Youden's index changes.

316  With the direct method, the choice of the prior has an influence on the shape of the posterior

317  probability, so that the threshold on $P(H_{ss}|d)$ can possibly be varied in a different interval (see

318  Figure 3 where $P(H_{ss}|d)$ never reaches 0 for example). However, in practice, there is no

319  influence on AUC, sensitivity and specificity because, again, the AUC depends essentially on

320  the distance distributions under $H_{ss}$ and $H_{ds}$.  The only palpable change is on the threshold

321  yielding the best compromise between sensitivity and specificity, threshold which adjusts to

322  $P(H_{ss})$ by roughly following it.

323  Thus, the assumption of equal prior probabilities has practically no impact on the LR estimate.

### 3.4. Limitations

325  From a practical point of view, our work suffers several limitations for a real-world forensic

326  application. First, for practical reasons, the subjects were sampled at a single time point, so

327    that the variability of the data is essentially due to the analytical variability. Second, the

328    chromatograms being compared are of the same nature, i.e. obtained on samples provided

329    by directly sampling the subjects (with contact with the subjects' hands) whereas in real life,

330    the unknown source sample will be obtained indirectly from an object on the crime scene

331    (without contact with the subject). Third, the odor collected on the crime scene might be

332    contaminated by other odors, from the environment or from other people. A study focused

333    on mixtures of odors, contaminations, and weathered traces has not been carried out yet but

334    is considered. However, despite these controlled conditions, the PCA of Figure 1 and the

335    statistics of Table 2 show that the data is already of limited reproducibility, so that the good

336    results we have obtained are encouraging concerning the robustness of the best method to

337    more realistic sampling conditions.

338    From a methodological point of view, the proposed methods are based on a common source

339    scenario, where it is asked whether the two traces originate from the same source or from

340    different sources without specifying which sources are considered, and not on a specific

341    source scenario, where the question is whether the two traces stem specifically from the

342    known source [Neuman and Ausdemore 2020]. The problem of the common source scenario

343    is that it does not take account of the typicality of the source, contrary to recommendations

344    for a better estimation of the strength of evidence through the LR [Morrison 2013, Tang and

345    Srihari 2014, Morrison and Enzinger 2018]. But to implement a specific source scenario, a

346    number of traces from the known source are needed in order to be able to estimate the

347    distribution under $H_{ss}$ (for the direct method) or to discriminate between the $H_{ss}$ and $H_{ds}$

348    populations (for the indirect methods), which is quite unpractical when dealing with human

349    hand odor, and not feasible at this stage of the study (at most four usable odor traces were

350    obtained, for only 253 subjects among the 534).

## 4. CONCLUSIONS

352    To summarize, the advantages expected from an indirect method are fully obtained, in

353    particular the dispensation to parameterize the likelihoods, and the robustness with respect

354    to differences in their variance and/or to possible outliers. Moreover, an increase in

355    performance of the indirect method as compared with the direct one is not obtained with a

356    scalar distance between odor traces, but when using the vector of the distances between each

357    feature of the odor traces. This improvement was not really expected, because, especially in

358 the forensic context, it is often advocated to convert multivariate data to a univariate datum

359 summarizing the relationship between features. Finally, the indirect method with a vectorial

360 distance proves also robust with respect to potentially irrelevant features since removing

361 them does not modify the performance, an appealing quality for dealing with traces which are

362 not yet solidly characterized, such as odor traces.

363 **REFERENCES**

364 T. Ali, L. Spreeuwers, R. Veldhuis, D. Meuwly (2015), Sampling variability in forensic likelihood-

365 ratio computation: A simulation study, Science & Justice 55 (6), 499-508,

366 https://doi.org/10.1016/j.scijus.2015.05.003

367 C. M. Bishop, Pattern Recognition and Machine Learning, Springer, 2006.

368 J. Chen, P. Zhang, M. Lv, H. Guo, Y. Huang, Z. Zhang, and F. Xu (2017), Influences of

369 Normalization Method on Biomarker Discovery in Gas Chromatography–Mass

370 Spectrometry-Based Untargeted Metabolomics: What Should Be Considered? Analytical

371 Chemistry 89 (10), 5342-5348 , https://doi.org/10.1021/acs.analchem.6b05152

372 C. Aitken, F. Taroni, Statistics and the Evaluation of Evidence for Forensic Scientists, John Wiley

373 & Sons, Second Edition (2004).

374 A. M. Curran, P. A. Prada, K. G. Furton (2010), The differentiation of the volatile organic

375 signatures of individuals through SPME-GC/MS of characteristic human scent compounds.

376 Journal of forensic sciences 55(1), 50–57, doi:10.1111/j.1556-4029.2009.01236.x

377 V. Cuzuel (2017), Développement d'une stratégie de caractérisation chimique de la signature

378 odorante d'individus par l'analyse chimiométrique de données issues de méthodes

379 séparatives bidimensionnelles, Thèse de doctorat de l'Université Pierre et Marie Curie,

380 ED388. https://hal.archives-ouvertes.fr/tel-01667972.

381 V. Cuzuel, G. Cognon, I. Rivals, C. Sauleau, F. Heulard, D. Thiébaut, J. Vial (2017), Origin,

382 analytical characterization and use of human odor in forensics, Journal of forensic sciences

383 62, 330–350, https://doi.org/10.1111/1556-4029.13394

384 V. Cuzuel, E. Portas, G. Cognon, I. Rivals, F. Heulard, D. Thiébaut, J. Vial (2017bis), Sampling

385 method development and optimization in view of human hand odor analysis by thermal

386 desorption coupled with gas chromatography and mass spectrometry., Analytical and

387 Bioanalytical Chemistry 409, 5113–5124. doi:10.1007/s00216-017-0458-8.

388 V. Cuzuel, G. Cognon, I. Rivals, F. Heulard, D. Thiébaut, J. Vial (2018), Human odor and
389    forensics. Optimization of a comprehensive gas chromatography method based on
390    orthogonality: how not to choose between criteria., J. Chromatogr. A. 1536 58-66,
391    doi:10.1016/j.chroma.2017.08.060.

392 V. Cuzuel, R. Leconte , G. Cognon, D. Thiébaut, J. Vial , C. Sauleau , I. Rivals (2018bis), Human
393    odor and forensics: towards Bayesian suspect identification using GCxGC-MS
394    characterization of hand odor, *Journal of Chromatography B* 1092, 379-385,
395    https://www.sciencedirect.com/science/article/abs/pii/S1570023218306068?via%3Dihub

396 H. E. Daniels (1944), The relation between measures of correlation in the universe of sample
397    per-mutation, Biometrika, Volume 33, Issue 2, 129-135,
398    https://doi.org/10.1093/biomet/33.2.129.

399 E. Enzinger, G. S. Morrison, F. Ochoa (2016) A demonstration of the application of the new
400    paradigm for the evaluation of forensic evidence under conditions reflecting those of a real
401    forensic-voice-comparison case, Science & Justice 56(1), 42-57,
402    https://doi.org/10.1016/j.scijus.2015.06.005.

403 H. Jeffreys, Theory of probability, Oxford University Press, Oxford, 1939.

404 T. Hastie, R. Tibshirani, Friedman J. The elements of statistical learning: data mining, inference,
405    and prediction, 2nd ed., Springer, New York, 2009.

406 A. J. Hanley, J.B. McNeil (1982), The Meaning and Use of the Area under a Receiver Operating
407    Characteristic (ROC) Curve, Radiology. 143, 29–36. doi:10.1148/radiology.143.1.7063747

408 Erwin J A T Mattijssen , Cilia L M Witteman, Charles E H Berger, Nicolaas W Brand, Reinoud D
409    Stoel (2020), Validity and Reliability of Forensic Firearm Examiners, Forensic Science
410    International 307 110112, https://doi.org/10.1016/j.forsciint.2019.110112

411 G. S. Morrison (2011), A comparison of procedures for the calculation of forensic likelihood
412    ratios from acoustic–phonetic data: Multivariate kernel density (MVKD) versus Gaussian
413    mixture model–universal background model (GMM–UBM), Speech Communication 53 (2),
414    242-256, https://doi.org/10.1016/j.specom.2010.09.005

415 G. S. Morrison (2013), Tutorial on logistic-regression calibration and fusion: converting a score
416    to a likelihood ratio. Australian Journal of Forensic Sciences, 45(2), 173-197.
417    https://doi.org/10.1080/00450618.2012.733025

418    G. S. Morrison, E. Enzinger (2018) Score based procedures for the calculation of forensic

419    likelihood ratios – Scores should take account of both similarity and typicality, Science &

420    Justice 58 (1), 47-58, https://doi.org/10.1016/j.scijus.2017.06.005.

421    C. Muehlethaler, G. Massonnet, T. Hicks (2016), Evaluation of infrared spectra analyses using

422    a likelihood ratio approach: A practical example of spray paint examination, Science &

423    Justice, 56 (2), 61-72, https://doi.org/10.1016/j.scijus.2015.12.001.

424    C. Neuman, C. Champod, R. Puch - Solis, N. Egli, A. Anthonioz, A. B. - Griffiths (2007),

425    Computation of likelihood ratios in fingerprint identification for configurations of any

426    number of minutiae, Journal of Forensic Sciences, 52(1), 54-64,

427    https://doi.org/10.1111/j.1556-4029.2006.00327.x

428    C. Neumann, Ian W Evett, J. Skerrett (2012), Quantifying the weight of evidence from a

429    forensic fingerprint comparison: A new paradigm, Journal of the Royal Statistical Society

430    Series A (Statistics in Society) 175(2), 1-26 https://doi.org/10.1111/j.1467-

431    985X.2011.01027.x

432    C. Neuman, M.A. Ausdemore (2020), Defence Against the Modern Arts: the Curse of Statistics

433    -Part II: "Score-based likelihood ratios", *Law, Probability and Risk* 19(1), 21–

434    42, https://doi.org/10.1093/lpr/mgaa006

435    F. Riva, C. Champod (2014), Automatic comparison and evaluation of impressions left by a

436    firearm on fired cartridge cases, Journal of forensic sciences 59(3), 637-47,

437    https://doi.org/10.1111/1556-4029.12382

438    F. Riva, E. J.A.T. Mattijssen, R. Hermsen, P. Pieper, W. Kerkhoff, C. Champod (2020),

439    Comparison and interpretation of impressed marks left by a firearm on cartridge cases –

440    Towards an operational implementation of a likelihood ratio based technique, Forensic

441    Science International 313, 110363, https://doi.org/10.1016/j.forsciint.2020.110363.

442    C. Robert, The Bayesian Choice: from Decision-Theoretic Motivations to Computational

443    Implementation, Springer, New York, 2001.

444    S.N. Srihari, C. Huang, H. Srinivasan (2008) On the Discriminability of the Handwriting of Twins,

445    Journal of Forensic Sciences 53(2), 430–446, https://doi.org/10.1111/j.1556-

446    4029.2008.00682.x

447   Y. Tang, S. N. Srihari (2014) Likelihood ratio estimation in forensic identification using similarity

448       and rarity, Pattern Recognition 47, 945-958, https://doi.org/10.1016/j.patcog.2013.07.014

449   Y. Tang , S. N. Shrihari (2014bis) Computational methods for the analysis of footwear

450       impression evidence, In "Computational Intelligence in Digital Forensics: Forensic

451       Investigation and Application", A. Muda, Y-H Choo, A. Abraham and S. N. Srihari (eds.),

452       Springer 2014, pages 333-383.

453   P. Vergeer, J. N. Hendrikse, M. M. P. Grutters, L. J. C. Peschier (2020). A method for forensic

454       gasoline comparison in fire debris samples: A numerical likelihood ratio system. Sci Justice.

455       2020 Sep;60(5):438-450, https://doi.org/10.1016/j.scijus.2020.06.002

456   W. J. Youden (1950). Index for rating diagnostic tests. Cancer. 3, 32–35,

457       https://doi.org/10.1002/1097-0142(1950)3:1<32::AID-CNCR2820030106>3.0.CO;2-3

458

459

460    **FIGURE CAPTIONS**

461

462    **Figure 1.** PCA of the data set showing the distribution of the odor traces between calibration

463    and test sets. The PCA was performed on the covariance matrix of the continuous features,

464    i.e. normalized in logarithmic scale.

465
466

467    **Figure 2**. Results of the direct method with the scalar distance d on binarized feature, as

468    functions of d, on the calibration set. a) Empirical distribution of d under $H_{ss}$ (1 594 pairs), and

469    estimated density (mixture of 2 Gaussians). b) Empirical distribution of d under $H_{ds}$

470    (841 457pairs), and estimated density (mixture of 2 Gaussians). c) Posterior probability of $H_{ss}$

471    obtained using Bayes' formula with equal priors. d) Likelihood ratio.

472
473

474    **Figure 3**. Results of the indirect method with the scalar distance d on continuous features, as

475    functions of d, on the calibration set. a) Empirical distribution of d under $H_{ss}$ (1 594 pairs). b)

476    Empirical distribution of d under $H_{ds}$ (841 457pairs). c) Logistic regression (dotted line), and

477    deduced posterior probability of $H_{ss}$ with equal priors (continuous line). d) "Likelihood ratio"

478    $\exp(a^T d + b)$ corresponding to the logistic regression (dotted line), and likelihood ratio

479    (continuous line).

480
481

482    **Figure 4**. Results of the indirect method with the vectorial distance d on continuous features,

483    as functions of the score $a^T d + b$, on the calibration set. a) Empirical distribution of the score

484    $a^T d + b$ under $H_{ss}$ (1 594 pairs). b) Empirical distribution of the score $a^T d + b$ under $H_{ds}$

485    (841 457pairs). c) Logistic regression (dotted line), and deduced posterior probability of $H_{ss}$

486    with equal priors (continuous line). d) "Likelihood ratio" $\exp(a^T d + b)$ corresponding to the

487    logistic regression (dotted line), and likelihood ratio (continuous line).

488
489

**TABLE CAPTIONS**

**Table 1.** Panel composition in terms of sex and age.

**Table 2.** Reproducibility of the odor trace features (peak areas), without and with two normalizations, estimated on the 490 subjects sampled at least twice (IQI stands for interquartile interval).

**Table 3.** Baseline comparison between the three methods on the calibration and test sets, using binarized features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's index, all in %.

**Table 4**. Baseline comparison between the three methods on the calibration and test sets, using continuous features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's index, all in %.

**Table 5**. Comparison between the three methods on the calibration and test sets, using binarized features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's index, with selection of the number of features among the 741 by cross-validation on the calibration set (AUC-CV3 denotes the mean 3-fold cross-validation AUC on the calibration set, and #feat. the number of features selected by cross-validation).

**Table 6.** Comparison between the three methods on the calibration and test sets, using continuous features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's index, with selection of the number of features among the 741 by cross-validation on the calibration set (AUC-CV3 denotes the mean 3-fold cross-validation AUC on the calibration set, and #feat. the number of features selected by cross-validation).

518 **TABLES**

519

520 **Table 1.** Panel composition in terms of sex and age.

521

| age\sex | 7-17 | 18-64 | 65-94 | total |
|---|---|---|---|---|
| man | 18 | 182 | 18 | 218 |
| woman | 28 | 268 | 20 | 316 |
| total | 46 | 450 | 38 | 534 |

522

523

524 **Table 2.** Repeatability of the odor trace features (peak areas), without and with two

525 normalizations, estimated on the 490 subjects sampled at least twice (IQI stands for

526 interquartile interval).

527

| Normalization | Median relative standard deviation [IQI] in % |
|---|---|
| None | 61.0 [32.1 ; 78.0] |
| In scalar scale | 56.1 [31.9 ; 74.2] |
| In logarithmic scale | 33.2 [17.9 ; 48.8] |

528

529

530

531

532

533

534

535

536

537

538

539

540 **Table 3.** Baseline comparison between the three methods on the calibration and test sets,

541 using binarized features, in terms of AUC, sensitivity (Sn) and specificity (Sp) for the threshold

542 maximizing Youden's index, all in %.

543

| Method | Calibration | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| | AUC | threshold | Sn | Sp | AUC | threshold | Sn | Sp |
| Direct | 91.2 | 0.43 | 80.6 | 91.3 | 91.4 | 0.54 | 78.4 | 94.9 |
| Indirect scal. d | 91.2 | 0.62 | 80.6 | 91.3 | 91.5 | 0.75 | 78.4 | 94.9 |
| Indirect vect. d | 98.5 | 0.54 | 93.4 | 96.3 | 97.1 | 0.56 | 91.1 | 94.2 |

544

545

546 **Table 4**. Baseline comparison between the three methods on the calibration and test sets,

547 using continuous features, in terms of AUC, sensitivity (Sn) and specificity (Sp) for the

548 threshold maximizing Youden's index, all in %.

549

| Method | Calibration | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| | AUC | threshold | Sn | Sp | AUC | threshold | Sn | Sp |
| Direct | 92.1 | 0.44 | 81.1 | 92.5 | 93.0 | 0.50 | 81.3 | 94.6 |
| Indirect scal. d | 92.1 | 0.64 | 81.1 | 92.5 | 93.0 | 0.73 | 81.3 | 94.6 |
| Indirect vect. d | 98.9 | 0.58 | 94.4 | 97.0 | 97.8 | 0.57 | 91.3 | 95.1 |

550

551

552

553

554

555

556

557

558 **Table 5**. Comparison between the three methods on the calibration and test sets, using
559 binarized features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's
560 index, with selection of the number of features among the 741 by cross-validation of the
561 calibration set (AUC-CV3 the mean 3-fold cross-validation AUC on the calibration set, and
562 #feat. denotes the number of features selected by cross-validation).

563

| Method | AUC-CV3 | #feat. | Calibration | | | Test | | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC | Sn | Sp | AUC | Sn | Sp |
| Direct | 94.4 | 267 | 94.5 | 80.9 | 94.7 | 94.4 | 84.6 | 92.8 |
| Indirect scal. d | 94.4 | 267 | 94.5 | 80.9 | 94.7 | 94.4 | 84.6 | 92.8 |
| Indirect vect. d | 96.4 | 440 | 97.6 | 89.5 | 97.5 | 97.0 | 91.3 | 94.3 |

564

565

566 **Table 6.** Comparison between the three methods on the calibration and test sets, using
567 continuous features, in terms of AUC, sensitivity (Sn) and specificity (Sp) maximizing Youden's
568 index, with selection of the number of features among the 741 by cross-validation of the
569 calibration set (AUC-CV3 the mean 3-fold cross-validation AUC on the calibration set, and
570 #feat. denotes the number of features selected by cross-validation).

571

| Method | AUC-CV3 | #feat. | Calibration | | | Test | | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC | Sn | Sp | AUC | Sn | Sp |
| Direct | 93.1 | 535 | 93.1 | 81.3 | 93.3 | 93.5 | 82.3 | 95.5 |
| Indirect scal. d | 93.1 | 535 | 93.1 | 81.3 | 93.3 | 93.5 | 82.3 | 95.5 |
| Indirect vect. d | 97.0 | 500 | 98.3 | 91.6 | 97.3 | 97.7 | 91.7 | 94.9 |

572
573

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
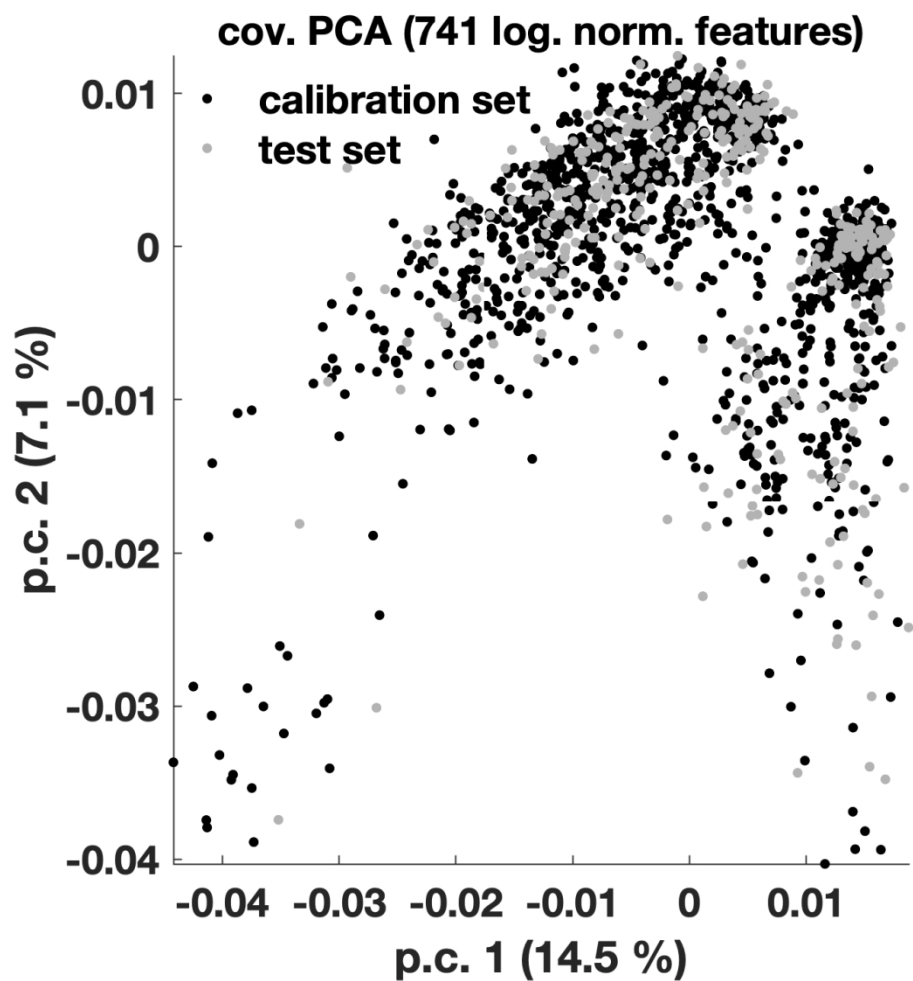51
52
53
54
55
56
57
58
59
60



Figure 1: PCA of the data set showing the distribution of the odor traces between calibration and test sets. The PCA was performed on the covariance matrix of the continuous features, i.e. normalized in logarithmic scale.
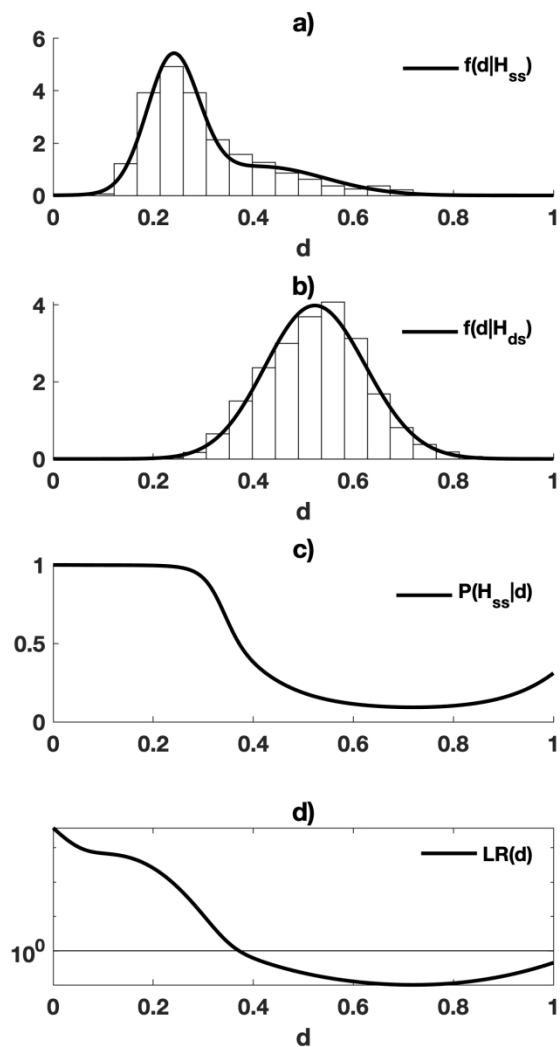
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Figure 2: Results of the direct method with the scalar distance d on binarized feature, as functions of d, on the calibration set. a) Empirical distribution of d under Hss (1 594 pairs), and estimated density (mixture of 2 Gaussians). b) Empirical distribution of d under Hds (841 457pairs), and estimated density (mixture of 2 Gaussians). c) Posterior probability of Hss obtained using Bayes' formula with equal priors. d) Likelihood ratio.
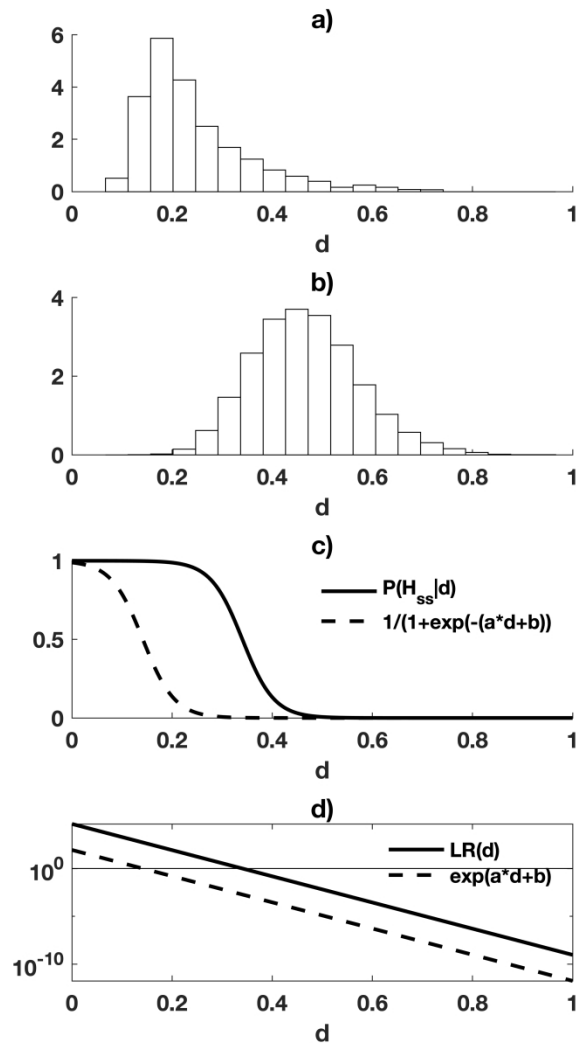
Figure 3: Results of the indirect method with the scalar distance d on continuous features, as functions of d, on the calibration set. a) Empirical distribution of d under Hss (1 594 pairs). b) Empirical distribution of d under Hds (841 457pairs). c) Logistic regression (dotted line), and deduced posterior probability of Hss with equal priors (continuous line). d) "Likelihood ratio" exp(aT d +b) corresponding to the logistic regression (dotted line), and likelihood ratio (continuous line).

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
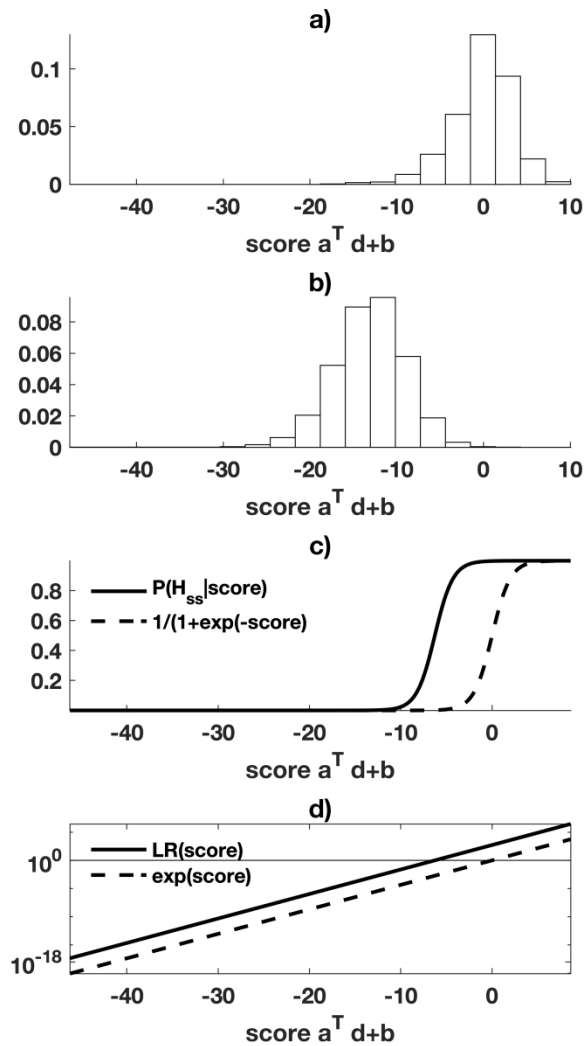24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

Figure 4: Results of the indirect method with the vectorial distance d on continuous features, as functions of the score aT d +b, on the calibration set. a) Empirical distribution of the score aT d + b under Hss (1 594 pairs). b) Empirical distribution of the score aT d + b under Hds (841 457pairs). c) Logistic regression (dotted line), and deduced posterior probability of Hss with equal priors (continuous line). d) "Likelihood ratio" exp(aT d +b) corresponding to the logistic regression (dotted line), and likelihood ratio (continuous line).

45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

| age \ sex | 7-17 | 18-64 | 65-94 | total |
|-----------|------|-------|-------|-------|
| man | 18 | 182 | 18 | 218 |
| woman | 28 | 268 | 20 | 316 |
| total | 46 | 450 | 38 | 534 |

| Normalization | Median relative standard deviation [IQI] in % |
|---|---|
| None | 61.0 [32.1 ; 78.0] |
| In scalar scale | 56.1 [31.9 ; 74.2] |
| In logarithmic scale | 33.2 [17.9 ; 48.8] |

| Method | Calibration | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| | AUC | threshold | Sn | Sp | AUC | threshold | Sn | Sp |
| Direct | 91.2 | 0.43 | 80.6 | 91.3 | 91.4 | 0.54 | 78.4 | 94.9 |
| Indirect scal. d | 91.2 | 0.62 | 80.6 | 91.3 | 91.5 | 0.75 | 78.4 | 94.9 |
| Indirect vect. d | 98.5 | 0.54 | 93.4 | 96.3 | 97.1 | 0.56 | 91.1 | 94.2 |

| | Calibration | | | | Test | | | |
|---|---|---|---|---|---|---|---|---|
| Method | AUC | threshold | Sn | Sp | AUC | threshold | Sn | Sp |
| Direct | 92.1 | 0.44 | 81.1 | 92.5 | 93.0 | 0.50 | 81.3 | 94.6 |
| Indirect scal. d | 92.1 | 0.64 | 81.1 | 92.5 | 93.0 | 0.73 | 81.3 | 94.6 |
| Indirect vect. d | 98.9 | 0.58 | 94.4 | 97.0 | 97.8 | 0.57 | 91.3 | 95.1 |

| Method | AUC-CV3 | #feat. | Calibration | | | Test | | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC | Sn | Sp | AUC | Sn | Sp |
| Direct | 94.4 | 267 | 94.5 | 80.9 | 94.7 | 94.4 | 84.6 | 92.8 |
| Indirect scal. d | 94.4 | 267 | 94.5 | 80.9 | 94.7 | 94.4 | 84.6 | 92.8 |
| Indirect vect. d | 96.4 | 440 | 97.6 | 89.5 | 97.5 | 97.0 | 91.3 | 94.3 |

| Method | AUC-CV3 | #feat. | Calibration | | | Test | | |
|---|---|---|---|---|---|---|---|---|
| | | | AUC | Sn | Sp | AUC | Sn | Sp |
| Direct | 93.1 | 535 | 93.1 | 81.3 | 93.3 | 93.5 | 82.3 | 95.5 |
| Indirect scal. d | 93.1 | 535 | 93.1 | 81.3 | 93.3 | 93.5 | 82.3 | 95.5 |
| Indirect vect. d | 97.0 | 500 | 98.3 | 91.6 | 97.3 | 97.7 | 91.7 | 94.9 |