



A Probabilistic Distance Clustering Algorithm Using Gaussian and Student-*t* Multivariate Density Distributions

Cristina Tortora¹ · Paul D. McNicholas² · Francesco Palumbo³

Received: 4 November 2019 / Accepted: 9 January 2020
© Springer Nature Singapore Pte Ltd 2020

Abstract

A new dissimilarity measure for cluster analysis is presented and used in the context of probabilistic distance (PD) clustering. The basic assumption of PD-clustering is that for each unit, the product between the probability of the unit belonging to a cluster and the distance between the unit and the cluster is constant. This constant is a measure of the classifiability of the point, and the sum of the constant over units is called joint distance function (JDF). The parameters that minimize the JDF maximize the classifiability of the units. The new dissimilarity measure is based on the use of symmetric density functions and allows the method to find clusters characterized by different variances and correlation among variables. The multivariate Gaussian and the multivariate Student-*t* distributions have been used, outperforming classical PD clustering, and its variation PD clustering adjusted for cluster size, on simulated and real datasets.

Keywords Cluster analysis · PD-clustering · Multivariate distributions · Dissimilarity measures

Introduction

Cluster analysis refers to a wide range of numerical methods aiming to find distinct groups of homogeneous units. Clustering in two or three dimensions is a natural task that humans can often do visually; however, machine approaches are needed for all but such low dimensions. We focus on partitioning clustering methods; given a number of clusters K , partitioning methods assign units to the K clusters optimizing a given criterion. These methods are generally divided into not model-based and model-based, according to the distributional assumptions. Model-based clustering or finite mixture model clustering assumes that the population probability density function is a convex linear combination of a finite number of density functions; accordingly, they are very well suited to clustering problems. A variety of methods and algorithms have been proposed for finite mixture

model parameter estimation. The most widely used strategy is to find the parameters that maximize the complete-data likelihood function using the expectation-maximization (EM) algorithm, which was proposed in 1977 [10] building on prior work (e.g., [5, 7, 23, 24]). The non-model-based methods generally optimize a criterion based on distance or dissimilarity measures. Different dissimilarity measures can be used based on the type of data, in this paper we focus on continuous data.

Formally, let us consider an $n \times J$ data matrix \mathbf{X} , with generic row vector $\mathbf{x}_i = (x_{i1}, \dots, x_{iJ})$. Partitioning algorithms aim to find a set of K clusters, \mathcal{C}_k , with $k = 1, \dots, K$, such that the elements inside a cluster are homogeneous and $\mathcal{C}_1 \cup \mathcal{C}_2 \cup \dots \cup \mathcal{C}_K = X$. If, for any pair $\{k, k'\} \in 1, \dots, K$, $\mathcal{C}_k \cap \mathcal{C}_{k'} = \emptyset$, then the clustering technique is called hard or crisp, otherwise it is called fuzzy or soft. In the latter case, each unit can belong to more than one cluster with certain membership degree.

The most frequently used non-model-based methods for continuous data are k-means [20] and its fuzzy analogue c-means [4], which minimize the sum of the within groups sum of squares over all variables. In spite of their simplicity, the optimal solution can only be found applying an iterative intuitively reasonable procedure. More recently, [3] proposed probabilistic distance (PD) clustering, a distribution free fuzzy clustering technique (i.e., non-model-based),

✉ Cristina Tortora
cristina.tortora@sjsu.edu

¹ Department of Mathematics and Statistics, San José State University, San José, CA, USA

² Department of Mathematics and Statistics, McMaster University, Hamilton, ON, Canada

³ Dipartimento di Scienze Politiche, University of Naples Federico II, Naples, Italy

59 where the membership degree is defined as heuristic proba- 107
 60 bility. PD clustering optimization problems represents a spe- 108
 61 cial case of the Weber–Fermat’s problem, when the number 109
 62 of the ‘attraction points’ is greater or equal to three, see [16] 110
 63 among others. In this framework, PD clustering assumes 111
 64 that the product of the probability of a point belonging to 112
 65 a cluster and the distance of the point from the center of 113
 66 the cluster is constant, and this constant is a measure of the 114
 67 classificability of the point. The method obtains the centers 115
 68 that maximize the classificability of all the points. A newer 116
 69 version of the algorithm that considers clusters of different 117
 70 size, PDQ-clustering, was proposed by [14] and an extension 118
 71 for high-dimensional data was proposed by [35, 36]. 119

72 Generally, non-model-based clustering techniques are only 120
 73 based on the distances between the points and the centers; 121
 74 therefore, they do not take into account the shape and the size of 122
 75 the clusters. Accordingly, these techniques may fail when clus- 123
 76 ters are either non-spherical or spherical with different radii. To 124
 77 overcome this issue we propose a new dissimilarity measure 125
 78 based on symmetric density functions that have the advantage 126
 79 of considering the variability and the correlation among the 127
 80 variables. We use two different density functions, the multi- 128
 81 variate Gaussian and the multivariate Student-*t*, but it could 129
 82 be extended to other symmetric densities. We then integrate 130
 83 this measure with PD-clustering and obtain new more flexible 131
 84 clustering techniques. Preliminary results can be found in [29]. 132

85 After a background section on PD-clustering and PDQ- 133
 86 clustering, Sect. "Background", we introduce the new dis- 134
 87 similarity measure and the new techniques, Sect. "Flexible 135
 88 Extensions of PD-Clustering". We then compare them with 136
 89 some model-based and distance-based algorithms on simu- 137
 90 lated and real datasets, Sect. Empirical Evidence from Simu- 138
 91 lated and Real Data. 139

92 **Background**

93 In this section we briefly introduce PD-clustering [3], a 140
 94 distance-based soft clustering algorithm, and its extension, 141
 95 PD-clustering adjusted for cluster size [14]. 142

96 **Probabilistic Distance Clustering**

97 Ben-Israel and Iyigun [3] proposed a non-hierarchical dis- 143
 98 tance-based clustering method, called probabilistic distance 144
 99 (PD) clustering. They then extended the method to account 145
 100 for clusters of different size, i.e., PDQ [14]. Tortora et al. 146
 101 [35] proposed a factor version of the method to deal with 147
 102 high-dimensional data. Recently, [29] further extended the 148
 103 method to include more flexibility. 149

104 In PD-clustering, the number of clusters *K* is assumed to 150
 105 be a priori known, and a wide review on how to choose *K* can 151
 106 be found in [8]. Given some random centers, the probability

of any point belonging to a cluster is assumed to be inversely 107
 proportional to the distance from the center of that cluster [13]. 108
 Suppose we have a data matrix **X** with *N* units and *J* variables, 109
 and consider *K* (non-empty) clusters. PD-clustering is based 110
 on two quantities: the distance of each data point **x_i** from each 111
 cluster centre **c_k**, denoted *d(x_i, c_k)*, and the probability of each 112
 point belonging to a cluster, i.e., *p(x_i, c_k)*, for *k* = 1, ..., *K* and 113
i = 1, ..., *N*. 114

For convenience, define *p_{ik}* := *p(x_i, c_k)* and *d_{ik}* := *d(x_i, c_k)*. 115
 PD-clustering is based on the principle that the product of 116
 the distances and the probabilities is a constant depending 117
 only on **x_i** [3]. Denoting this constant as *F(x_i)*, we can write 118
 this principle as 119

$$p_{ik}d_{ik} = F(x_i), \tag{1}$$

where *F(x_i)* depends only on **x_i**, i.e., *F(x_i)* does not depend 122
 on the cluster *k*. As the distance from the cluster centre 123
 decreases, the probability of the point belonging to the 124
 cluster increases. The quantity *F(x_i)* is a measure of the 125
 closeness of **x_i** to the cluster centres, and it determines the 126
 classificability of the point **x_i** with respect to the centres **c_k**, 127
 for *k* = 1, ..., *K*. The smaller the *F(x_i)* value, the higher the 128
 probability of the point belonging to one cluster. If all of the 129
 distances between the point **x_i** and the centers of the clusters 130
 are equal to *d_i*, then *F(x_i)* = *d_i*/*K* and all of the probabilities 131
 of belonging to each cluster are equal, i.e., *p_{ik}* = 1/*K*. The 132
 sum of *F(x_i)* over *i* is called joint distance function (JDF). 133
 Starting from (1), it is possible to compute *p_{ik}*, i.e., 134

$$p_{ik} = \frac{\prod_{m \neq k} d_{im}}{\sum_{m=1}^K \prod_{r \neq m} d_{ir}}, \tag{2}$$

for *k* = 1, ..., *K*, and *i* = 1, ..., *N*. The whole clustering prob- 137
 lem consists in the identification of the centers that minimize 138
 the JDF: 139

$$JDF = \sum_{i=1}^n \sum_{k=1}^K d_{ik}p_{ik}. \tag{3}$$

Extensive details on PD clustering are given in [3], who 142
 suggest using *p*² in (3) because it is a smoothed version of 143
 the problem. It follows that the optimized functions become 144

$$JDF = \sum_{i=1}^n \sum_{k=1}^K d_{ik}p_{ik}^2, \tag{4}$$

and the centers can be computed as 147

$$c_k = \frac{\sum_{i=1}^N u_{ik}x_i}{\sum_{j=1}^N u_{jk}}, \tag{5}$$

with *u_{ik}* = *p_{ik}*²/*d_{ik}*. 149
 150

151 It is worth noting that the function p_{ik} respects all neces- 188
 152 sary conditions to be a probability and yet no assumptions 189
 153 are made on the distribution of this function; further, p_{ik} can 190
 154 only be computed given \mathbf{x}_i and for every \mathbf{c}_k [28]. Following 191
 155 [13], we refer to p_{ik} as subjective probabilities, which are 192
 156 based on degree of belief (see [2]).

157 **PD Clustering Adjusted for Cluster Size**

158 The probabilities obtained using 1 do not consider the clus- 195
 159 ter size, and the algorithm tends to fail when clusters are 196
 160 unbalanced. Moreover, the resulting clusters have similar 197
 161 variance and covariance matrices. To overcome these issues 198
 162 [14] proposed PD-clustering adjusted for cluster size (PDQ). 199
 163 They assume that

164
$$\frac{p_{ik}^2 d_{ik}}{q_k} = F(\mathbf{x}_i), \tag{6}$$

165 where q_k is the cluster size, with the constraint that 200
 166 $\sum_{k=1}^K q_k = N$. The p_{ik} can then be computed via 201

168
$$p_{ik} = \frac{\prod_{m \neq k} d_{im}/q_m}{\sum_{m=1}^K \prod_{r \neq m} d_{ir}/q_r}. \tag{7}$$

169 The cluster size is considered a variable, the value of q_k that 202
 170 minimizes (6) is 203

172
$$q_k = N \frac{\left(\sum_{i=1}^N d_{ik} p_{ik}^2\right)^{1/2}}{\sum_{k=1}^K \left(\sum_{i=1}^N d_{ik} p_{ik}^2\right)^{1/2}}, \tag{8}$$

173 for $k = 1, \dots, K - 1$, and 204

175
$$q_k = N - \sum_{k=1}^{K-1} q_k.$$

177 **Flexible Extensions of PD-Clustering**

178 **Gaussian PD-Clustering**

179 The PDQ algorithm can detect clusters of different size and 210
 180 with different within-cluster variability; however, it can still 211
 181 fail at detecting the clustering partition when variables are 212
 182 correlated or when there are outliers in the data. To over- 213
 183 come these issues we proposed a new dissimilarity measure 214
 184 based on a density function. Let $M_k = \max\{f(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\theta}_k)\}$ and 215
 185 define the quantity

186
$$\delta_{ik} = \log \left(M_k f(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\theta}_k)^{-1} \right), \tag{9}$$

188 which is a dissimilarity measure where $f(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\theta}_k)$ is a sym- 189
 190 metric unimodal density function with location parameter $\boldsymbol{\mu}_k$ 191
 and parameter vector $\boldsymbol{\theta}_k$. See appendix for the proof.

192 Recall that the density of a multivariate Gaussian distribu- 193
 194 tion is

193
$$\phi(\mathbf{x}_i; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{J/2}} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}, \tag{10}$$

195 and define $\phi_{ik} := \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$, where $k = 1, \dots, K$. Using 196
 197 (10) in (9) and the result in (6), the JDF becomes

197
$$\begin{aligned} \text{JDF} = & \sum_{i=1}^n \sum_{k=1}^K \frac{p_{ik}^2}{q_k} \log(M_k) + \sum_{i=1}^n \sum_{k=1}^K \frac{1}{2} \frac{p_{ik}^2}{q_k} \log((2\pi)^J |\boldsymbol{\Sigma}_k|) \\ & + \sum_{i=1}^n \sum_{k=1}^K \frac{1}{2} \frac{p_{ik}^2}{q_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k). \end{aligned} \tag{11}$$

199 This technique is called Gaussian PD-clustering, GPDC, and 200
 201 it offers many advantages when compared to PD-clustering. 202
 203 The new dissimilarity measure already takes into account 204
 the impact of different within cluster variances and the cor- 205
 relation among variables.

204 The clustering problem now consists in the estimation of $\boldsymbol{\mu}_k$ 205
 206 and $\boldsymbol{\Sigma}_k$, with $k = 1, \dots, K$, that minimize (11). A differentia- 206
 207 tion procedure leads to these estimates. An iterative algorithm 207
 208 is then used to compute the belonging probabilities and update 208
 the parameter estimates. More specifically, differentiating (11) 209
 with respect to $\boldsymbol{\mu}_k$ gives

210
$$\frac{\partial \text{JDF}}{\partial \boldsymbol{\mu}_k} = -\frac{1}{2} \sum_{i=1}^n \frac{p_{ik}^2}{q_k} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k). \tag{12}$$

211 Setting (12) equal to zero and solving for $\boldsymbol{\mu}_k$ gives 212

213
$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^n p_{ik}^2 \mathbf{x}_i}{\sum_{i=1}^n p_{ik}^2} \tag{13}$$

214 Now, differentiating (11) with respect to $\boldsymbol{\Sigma}_k$ gives 215

216
$$\begin{aligned} \frac{\partial \text{JDF}}{\partial \boldsymbol{\Sigma}_k} = & \sum_{i=1}^n \frac{1}{2} \frac{p_{ik}^2}{q_k} \boldsymbol{\Sigma}_k^{-1} - \boldsymbol{\Sigma}_k^{-1} \sum_{i=1}^n \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)' \frac{p_{ik}^2}{q_k} \boldsymbol{\Sigma}_k^{-1} \\ = & \frac{1}{2} \boldsymbol{\Sigma}_k^{-1} \left[\sum_{i=1}^n \frac{p_{ik}^2}{q_k} - \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)' p_{ik}^2 \boldsymbol{\Sigma}_k^{-1} \right]. \end{aligned} \tag{14}$$

217 Setting (14) equal to zero and solving for $\boldsymbol{\Sigma}_k$ gives 218

219
$$\boldsymbol{\Sigma}_k = \frac{\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)' p_{ik}^2}{\sum_{i=1}^n p_{ik}^2}. \tag{15}$$

220 It follows that, at generic iteration $(t + 1)$, the parameters that 221
 222 minimize the (11) are:

Author Proof

$$\mu_k^{(t+1)} = \frac{\sum_{i=1}^n p_{ik}^2 x_i}{\sum_{i=1}^n p_{ik}^2}, \quad (16)$$

$$\Sigma_k^{(t+1)} = \frac{\sum_{i=1}^n (x_i - \mu_k^{(t+1)})(x_i - \mu_k^{(t+1)})' p_{ik}^2}{\sum_{i=1}^n p_{ik}^2}. \quad (17)$$

Our iterative procedure for Gaussian mixture model-based clustering parameter estimation can be summarized as follows:

Algorithm 1 GPDC

```

1: procedure GPDC(X, K)
2:   for k = 1, ..., K do                                ▷ Initialization
3:     Random initialization of μk
4:     set Σk as identity matrix
5:   while μk(t) ≠ μk(t+1) do                            ▷ Core
6:     for k = 1, ..., K do
7:       update qk according to (8)
8:       update pik according to (7)
9:       update μk according to (16)
10:      update Σk according to (17)
11:   return pik, μk, Σk
    
```

Generalization to a Multivariate Student-t Distribution

The same procedure can be generalized to any symmetric distribution. In this subsection we use the multivariate Student-t distribution, generating an algorithm identified as Student-t PD-Clustering (TPDC). TPDC can detect clusters characterized by heavy tails; furthermore, the Student-t distribution has been often used on datasets characterized by outliers [17]. Now, replace (10) with a multivariate Student-t distribution, i.e.,

$$f(x, \mu, \Sigma, \nu) = \frac{\Gamma\left(\frac{\nu+J}{2}\right) |\Sigma|^{-\frac{1}{2}}}{(\pi\nu)^{\frac{1}{2}J} \Gamma\left(\frac{\nu}{2}\right) \left\{1 + \frac{\delta(x, \mu, \Sigma)}{\nu}\right\}^{\frac{1}{2}(\nu+J)}}, \quad (18)$$

where $\delta(x, \mu, \Sigma) = (x - \mu)' \Sigma^{-1} (x - \mu)$, and proceed as in Sect. 3.1 Then, the JDF becomes:

$$\begin{aligned} \text{JDF} = & \sum_{i=1}^n \sum_{k=1}^K \frac{p_{ik}^2}{q_k} \log(M_k) \\ & + \sum_{i=1}^n \sum_{k=1}^K \frac{p_{ik}^2}{q_k} \left[-\log \left\{ \Gamma\left(\frac{\nu_k+J}{2}\right) |\Sigma_k|^{-\frac{1}{2}} \right\} \right. \\ & \left. + \sum_{i=1}^n \sum_{k=1}^K \frac{p_{ik}^2}{q_k} \log \left\{ (\pi\nu_k)^{\frac{1}{2}J} \Gamma\left(\frac{\nu_k}{2}\right) \left(1 + \frac{\delta(x_i, \mu_k, \Sigma_k)}{\nu_k}\right)^{\frac{\nu_k+J}{2}} \right\} \right]. \end{aligned} \quad (19)$$

The parameters that optimize (19) can be found by differentiating with respect to μ_k , Σ_k , and ν_k , respectively. Specifically, at a generic iteration $(t + 1)$, the parameters that minimize (19) are:

$$\mu_k^{(t+1)} = \frac{\sum_{i=1}^n w_{ik} x_i}{\sum_{i=1}^n w_{ik}}, \quad (20)$$

with $w_{ik} = p_{ik}^2 / [v_k^{(t)} + \delta(x_i, \mu_k^{(t)}, \Sigma_k^{(t)})]$,

$$\Sigma_k^{(t+1)} = \frac{\sum_{i=1}^n p_{ik}^2 (x_i - \mu_k^{(t+1)})(x_i - \mu_k^{(t+1)})' s_{ik}}{\sum_{i=1}^n p_{ik}^2}, \quad (21)$$

with $s_{ik} = (v_k^{(t)} + J) / [v_k^{(t)} + \delta(x_i, \mu_k^{(t+1)}, \Sigma_k^{(t)})]$, and the degrees of freedom update $v_k^{(t+1)}$ is the solution to the following equation:

$$\begin{aligned} & \sum_{i=1}^n p_{ik}^2 \left[\psi\left(\frac{\nu_k}{2}\right) - \psi\left(\frac{\nu_k+J}{2}\right) + \frac{J}{2\nu_k} \right] \\ & + \sum_{i=1}^n p_{ik}^2 \left[\frac{1}{2} \log \left(1 + \frac{\delta(x_i, \mu_k^{(t+1)}, \Sigma_k^{(t+1)})}{\nu_k^{(t)}} \right) \right] \\ & - \frac{1}{2} \frac{\nu_k + J}{\nu_k} \sum_{i=1}^n p_{ik}^2 \frac{\delta(x_i, \mu_k^{(t+1)}, \Sigma_k^{(t+1)})}{\nu_k^{(t)} + \delta(x_i, \mu_k^{(t+1)}, \Sigma_k^{(t+1)})} = 0, \end{aligned} \quad (22)$$

where

$$\psi(\nu) = \left(\frac{1}{\Gamma(\nu)} \right) \frac{\delta \Gamma(\nu)}{\delta \nu}.$$

Our iterative algorithm can be summarized as follows:

Algorithm 2 TPDC

```

1: procedure TPDC(X, K)
2:   for k = 1, ..., K do                                ▷ Initialization
3:     Random initialization of μk
4:     set Σk as identity matrix
5:     νk = 20
6:   while μk(t) ≠ μk(t+1) do                            ▷ Core
7:     for k = 1, ..., K do
8:       update qk according to (8)
9:       update pik according to (7)
10:      update μk according to (20)
11:      update Σk according to (21)
12:      update νk solving (22)
13:   return pik, μk, Σk, νk
    
```

Algorithm Details

All the proposed techniques require a random initialization. Random starts can lead to unstable solutions, to avoid this

Author Proof

269 problem the algorithms use multiple starts. Moreover, the
 270 functions include the option to use PD-clustering or parti-
 271 tion around medoids (PAM; [15]) to start. As for many other
 272 clustering techniques, the optimized function, the JDF in
 273 (4), is not convex—not even quasi-convex—and may have
 274 other stationary points. For a fixed value of Σ_k , the JDF is a
 275 monotonically decreasing function, this guarantees that the
 276 function converges to a minimum, not necessarily a global
 277 minimum. The proposed techniques, GPDC and TPDC,
 278 introduce the estimate of Σ_k , giving much more flexibility,
 279 albeit the JDF is no longer monotonically decreasing. Using
 280 (9) in (4), we obtain

$$281 \text{JDF} = \sum_{i=1}^n \sum_{k=1}^K p_{ik}^2 (\log M_k - \log \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))$$

282 with $M_k \geq \phi(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$. Therefore, for every $k = 1, \dots, K$,
 283 the function is upper-bounded for non-degenerate density
 284 functions. The convergence of the algorithm cannot depend
 285 on the JDF but is based on $\boldsymbol{\mu}_k$. The time complexity of the
 286 algorithm is comparable to the EM algorithm, both algo-
 287 rithms require the inversion and the determinant of a $J \times J$
 288 matrix, therefore, the time complexity is of $O(n^3JK)$, where
 289 n is the number of observations, J the number of variables,
 290 and K the number of clusters.
 291

292 Empirical Evidence from Simulated and Real 293 Data

294 The proposed algorithm has been evaluated on real and simu-
 295 lated datasets. The simulated datasets have been used to
 296 illustrate the ability of the algorithms to recover the param-
 297 eters of the distributions and to compare the new techniques
 298 with some existing methods. In the following sessions we
 299 used the software R [26], the functions for both GPDC and
 300 TPDC are included in the R package `FPDclustering`
 301 [37].

302 Simulation Study

303 The same design was used twice, the first time each clus-
 304 ter was generated from a multivariate Gaussian distribution
 305 with three variables and $K = 3$ clusters. The second time,
 306 using a multivariate Student- t distribution with five degrees
 307 of freedom, same number of variables and clusters. We set
 308 the parameter using a four factor full factorial design. There
 309 are two factors per each level, where the levels are

- 310 – Overlapping and not overlapping clusters
- 311 – Different number of elements per clusters
- 312 – Unitary variance and variance bigger than 1
- 313 – Uncorrelated and correlated variables

Table 1 Model parameters used to generate the simulated datasets

Not overlapping	Overlapping
$\boldsymbol{\mu}_1 = (0, 0, 0)'$	$\boldsymbol{\mu}_1 = (0, 0, 0)'$
$\boldsymbol{\mu}_2 = (-7, 7, 0)'$	$\boldsymbol{\mu}_2 = (-4, 4, 0)'$
$\boldsymbol{\mu}_3 = (-7, 0, 7)'$	$\boldsymbol{\mu}_3 = (-4, 0, 4)'$
Option 1	Option 2
$n_1 = 200$	$n_1 = 200$
$n_2 = 300$	$n_2 = 100$
$n_3 = 100$	$n_3 = 300$
Unitary variance	bigger than 1
$\text{diag}(\boldsymbol{\Sigma}_1) = 1$	$\text{diag}(\boldsymbol{\Sigma}_1) = 1$
$\text{diag}(\boldsymbol{\Sigma}_2) = 1$	$\text{diag}(\boldsymbol{\Sigma}_2) = 16$
$\text{diag}(\boldsymbol{\Sigma}_3) = 1$	$\text{diag}(\boldsymbol{\Sigma}_3) = 2.25$
Not correlated	Correlated
$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\boldsymbol{\Sigma}_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$
$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\boldsymbol{\Sigma}_2 = \begin{pmatrix} 1 & -0.5 & -0.5 \\ -0.5 & 1 & 0.5 \\ -0.5 & 0.5 & 1 \end{pmatrix}$
$\boldsymbol{\Sigma}_3 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$	$\boldsymbol{\Sigma}_3 = \begin{pmatrix} 1 & 0.7 & 0.7 \\ 0.7 & 1 & 0.7 \\ 0.7 & 0.7 & 1 \end{pmatrix}$

Table 1 shows the parameters used in the simulation study. 314

The datasets have been generated using the R package `mvtnorm` [11]. Tables 5, 6, 7, 8, 9, 10, 11, 12 in Appendix B.2 show the true and the average estimated values of the parameters obtained from 50 runs of the GPDC and TPDC algorithms. For sake of space, comments are limited to groups of scenarios. The factors that affect the estimates the most are the change in variances and the amount of overlap. Specifically, when data are simulated using multivariate Gaussian distributions, in cases 5–8 and 13–16, the variances are not homogeneous and the GPDC tends to underestimate the bigger variances and overestimate the smaller ones. The TPDC is less affected by this issue, i.e., it underestimates some of the variances but the degrees of freedom recover; however, in the two extreme scenarios, 8 and 16, it cannot recover the cluster structures. Similar outcomes occur when data are simulated using a multivariate Student- t distribution; moreover, as expected on those datasets, the GPDC tends to overestimate the variances and TPDC tends to underestimate the variances and compensate with the degrees of freedom. 315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335

On the same datasets we used the functions `gpcm`, option `VVV`, of the R package `mixture` [6] for the Gaussian mixture models (GMM) and the function `teigen`, option `UUUU`, of the homonymous R package [1] for the mixtures of multivariate Student- t distributions (TMM). The k-means 336
337
338
339
340

Table 2 Average ARI and standard deviation on 50 datasets per scenario

	Correlated	Unitary Variance	Over-lapping	n_k	GPDC		TPDC		PDQ		k-means		GMM		TMM	
					ARI	SD	ARI	SD	ARI	SD	ARI	SD	ARI	SD	ARI	SD
Multivariate Gaussian distribution																
1	No	Yes	No	Op 1	1.00	0.00	1.00	0.00	1.00	0.00	0.91	0.20	0.99	0.06	1.00	0.00
2	No	Yes	no	Op 2	1.00	0.00	1.00	0.00	1.00	0.00	0.92	0.18	0.96	0.12	1.00	0.00
3	No	Yes	Yes	Op 1	0.97	0.01	0.98	0.01	0.98	0.01	0.99	0.01	0.99	0.01	0.99	0.01
4	No	Yes	Yes	Op 2	0.97	0.01	0.98	0.01	0.98	0.01	0.98	0.01	0.99	0.01	0.99	0.01
5	No	No	No	Op 1	0.80	0.03	0.84	0.03	0.74	0.04	0.61	0.08	0.93	0.07	0.93	0.02
6	No	No	No	Op 2	0.94	0.02	0.93	0.02	0.74	0.05	0.89	0.09	0.97	0.01	0.97	0.01
7	No	No	Yes	Op 1	0.59	0.04	0.60	0.05	0.30	0.04	0.35	0.06	0.76	0.08	0.76	0.04
8	No	No	Yes	Op 2	0.56	0.12	0.50	0.09	0.60	0.06	0.75	0.08	0.87	0.02	0.87	0.02
9	Yes	Yes	No	Op 1	1.00	0.00	1.00	0.00	1.00	0.00	0.86	0.22	0.95	0.14	1.00	0.00
10	Yes	Yes	No	Op 2	1.00	0.00	1.00	0.00	1.00	0.00	0.91	0.19	0.91	0.18	1.00	0.00
11	Yes	Yes	Yes	Op 1	0.97	0.02	0.98	0.02	0.98	0.01	0.93	0.15	0.96	0.09	0.99	0.01
12	Yes	yes	Yes	Op 2	0.98	0.02	0.99	0.01	0.98	0.01	0.98	0.07	0.98	0.08	1.00	0.00
13	Yes	No	No	Op 1	0.76	0.08	0.90	0.03	0.67	0.05	0.52	0.08	0.95	0.07	0.96	0.01
14	Yes	No	No	Op 2	0.97	0.02	0.96	0.02	0.77	0.06	0.91	0.02	0.99	0.01	0.99	0.01
15	Yes	No	Yes	Op 1	0.44	0.08	0.45	0.08	0.33	0.05	0.32	0.04	0.82	0.09	0.83	0.03
16	Yes	No	Yes	Op 2	0.54	0.07	0.45	0.08	0.54	0.08	0.80	0.10	0.94	0.01	0.94	0.01
Multivariate Student- <i>t</i> distribution, 5 df																
1	no	Yes	No	Op 1	0.98	0.01	0.99	0.01	0.99	0.01	0.98	0.07	0.97	0.06	0.99	0.01
2	No	Yes	No	Op 2	0.99	0.01	0.99	0.01	0.99	0.01	0.93	0.17	0.95	0.11	0.99	0.01
3	No	Yes	Yes	Op 1	0.88	0.02	0.89	0.02	0.91	0.02	0.90	0.02	0.90	0.02	0.91	0.02
4	No	Yes	Yes	Op 2	0.88	0.03	0.89	0.02	0.90	0.02	0.89	0.07	0.89	0.05	0.90	0.02
5	No	No	No	Op 1	0.92	0.02	0.94	0.02	0.94	0.02	0.89	0.07	0.93	0.07	0.95	0.01
6	No	No	No	Op 2	0.96	0.02	0.96	0.01	0.93	0.02	0.96	0.01	0.91	0.06	0.96	0.01
7	No	No	Yes	Op 1	0.69	0.04	0.71	0.04	0.70	0.04	0.63	0.04	0.56	0.11	0.76	0.03
8	No	No	Yes	Op 2	0.76	0.03	0.76	0.03	0.62	0.05	0.81	0.03	0.66	0.06	0.78	0.05
9	Yes	Yes	No	Op 1	0.99	0.01	0.99	0.01	0.99	0.01	0.90	0.19	0.92	0.17	0.99	0.01
10	Yes	Yes	No	Op 2	0.99	0.01	0.99	0.01	0.99	0.01	0.87	0.22	0.93	0.14	0.99	0.01
11	Yes	Yes	Yes	Op 1	0.87	0.03	0.88	0.03	0.90	0.02	0.83	0.16	0.90	0.07	0.93	0.02
12	Yes	Yes	Yes	Op 2	0.87	0.05	0.89	0.04	0.93	0.02	0.92	0.11	0.94	0.06	0.96	0.01
13	Yes	No	No	Op 1	0.93	0.02	0.95	0.02	0.92	0.02	0.79	0.17	0.93	0.10	0.96	0.01
14	Yes	No	No	Op 2	0.97	0.01	0.98	0.01	0.95	0.02	0.97	0.01	0.96	0.03	0.98	0.01
15	Yes	No	Yes	Op 1	0.56	0.10	0.59	0.08	0.65	0.04	0.55	0.08	0.72	0.13	0.82	0.03
16	Yes	No	Yes	Op 2	0.49	0.29	0.51	0.32	0.56	0.31	0.67	0.37	0.64	0.35	0.71	0.38

341 algorithm is part of the stats package [27], and the PDQ-
 342 clust function for PDQ clustering is part of the FPD-
 343 clustering package [37].

344 To compare the clustering performance of the methods
 345 we used the adjusted Rand index (ARI) [12]. It compares
 346 predicted classifications with true classes. The ARI corrects
 347 the Rand index [30] for chance, its expected value under ran-
 348 dom classification is 0, and it takes a value of 1 when there
 349 is perfect class agreement. Steinley [31] gives guidelines for
 350 interpreting ARI values. Table 2 shows the average ARI and
 351 the standard deviation on 50 runs for each algorithm.

352 As pointed out in the previous sections, GPDC and
 353 TPDC are framed in a non-parametric view; however, to

354 evaluate the performance we compare them with the GMM
 355 and TMM. The performance is not expected to be better
 356 than those techniques, although in most scenarios GPDC
 357 and TPDC perform as well as finite mixture models. As
 358 expected, k-means results are impacted by correlations and
 359 not homogeneous variances. PDQ cannot recover the correct
 360 clustering partition in case of overlapping and not homoge-
 361 neous variance. It is not affected by changes in group size
 362 or correlation. The proposed techniques GPDC and TPDC
 363 outperform k-means and PDQ in most scenarios, they show
 364 weakness in the two most extreme situations, i.e., scenarios
 365 8 and 16. Specifically, when clusters have different variances
 366 and the biggest variance is associated with the smallest

Table 3 Number of units, variables, and clusters for the three real datasets

	n	J	K	
Seed	210	4	3	UCI machine learning repository
HSCT	9702	3	4	Terry Fox Lab
AIS	202	4	2	EMMIXuskew R package

Table 4 Adjusted Rand index for the real datasets

	GPDC	TPDC	PDQ	k-means	GMM	TMM
Seed	0.41	0.33	0.17	0.16	0.44	0.53
HSCT	0.99	0.98	0.98	0.72	0.88	0.99
AIS	0.88	0.90	0.16	0.06	0.72	0.81

367 cluster, they fail detecting the clustering partitions. Fig-
 368 ures 5, 6, 7, 8, 9, 10, 11, 12 in Appendix B.2 show examples
 369 of simulated datasets for each scenario.

370 **Real Data Analysis**

371 We performed a real data analysis on three datasets that
 372 differ in size and number of clusters (details in Table 3).
 373 We performed variable selection prior to cluster analysis
 374 (details in Appendix B.1). The seed dataset¹ contains infor-
 375 mation about kernels belonging to three different varieties of
 376 wheat—Kama, Rosa and Canadian—with 70 observations
 377 per variety (see Fig. 1). We used the variables: compactness,
 378 length of kernel, width of kernel, and asymmetry coefficient.
 379 The hematopoietic stem cell transplant (HSCT) data were
 380 collected in the Terry Fox Lab at the British Columbia Can-
 381 cer Agency. It contains information about 9780 cells, each
 382 stained with four fluorescent dyes. Experts identified four
 383 clusters; moreover, 78 cells were deemed “dead”, leaving a
 384 total of 9702 observation, we selected the three most infor-
 385 mative variables. Figure 2 shows the partitions defined by the
 386 experts. The Australian Institute of sport dataset² contains
 387 data on 102 male and 100 female athletes for the Australian
 388 institute of sports. We selected the variables: height in cm,
 389 hematocrit, plasma ferritin concentration, and percent body
 390 fat, see Fig. 3.

391 Table 4 shows the ARI on the three datasets. On the seed
 392 dataset, GPDC and TPDC perform better than PDQ and
 393 k-means. The improvement from PDQ is noticeable, PDQ
 394 gives an ARI of 0.17, while GPDC gives an ARI of 0.41.
 395 On this dataset, TMM gives the best performance. On the
 396 HSCT dataset, GPDC, TPDC, PDQ, and TMM have a very

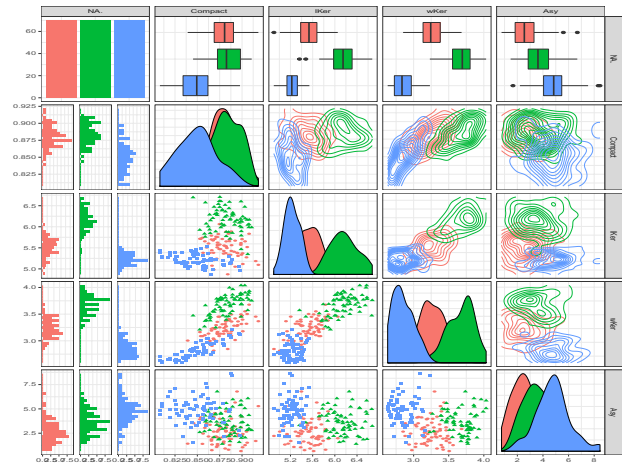


Fig. 1 Seed dataset, each color and symbol representing a different variety of wheat

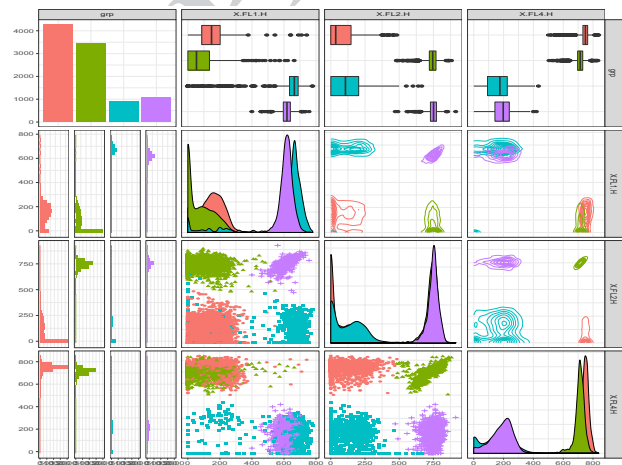


Fig. 2 HSCT dataset, each color and symbol representing a partition defined by the experts

high ARI. On the AIS dataset, GPDC and TPDC give the
 397 best performance. 398

399 **Conclusion**

A new distance measure based on density functions is intro-
 400 duced and used in the context of probabilistic distance clus-
 401 tering adjusted for cluster size (PDQ). PDQ assumes that,
 402 for a generic unit, the product between the probability of
 403 belonging to a cluster and its distance from the cluster is
 404 constant. The minimization of the sum of these constants
 405 over the units leads to clusters that maximize the classifi-
 406 ability of the data. We introduce two algorithms based on
 407 PDQ that use distance measures based on the multivariate
 408

1FL01 ¹ <http://archive.ics.uci.edu/ml/>.

2FL01 ² GLMsData R package.

Author Proof

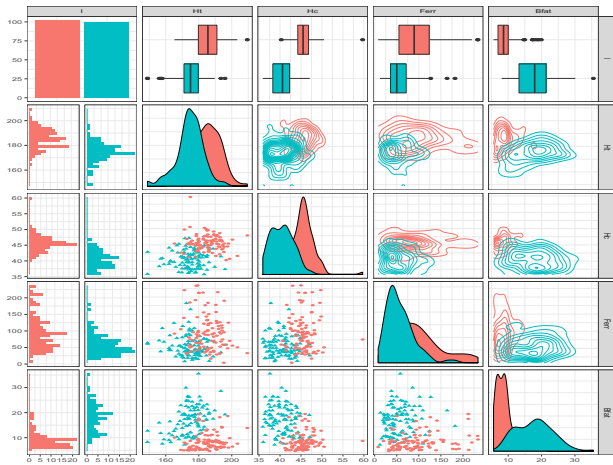


Fig. 3 AIS dataset, each color and symbol representing male and female athletes

409 Gaussian distribution and on the multivariate Student-*t* dis-
 410 tribution. Using simulated and real datasets we show how
 411 the new algorithms over-perform PDQ and the well known
 412 k-means algorithm.

413 The algorithm could be extended using different distribu-
 414 tions. Further to this point, we mentioned outliers as a possi-
 415 ble motivation for the PDQ approach with the multivariate
 416 Student-*t* distribution (Sect. 3). However, if the objective is
 417 dealing with outliers, it will be better to consider the PDQ
 418 approach with the multivariate contaminated normal dis-
 419 tribution [25] and this will be a topic of future work. Other

420 approaches for handling cluster concentration will also be
 421 considered (e.g., [9]) as will methods that accommodate
 422 asymmetric, or skewed, clusters (e.g., [18, 19, 21, 22, 32,
 423 34]).
 424

425 **Funding** this study was funded by a Discovery Grant from the Natural
 426 Sciences and Engineering Research Council of Canada and the Canada
 427 Research Chairs program (McNicholas). During the development of
 428 the present work, Prof. Francesco Palumbo had a short term visiting at
 429 the San Jose State University (CA) financially supported by the Inter-
 430 national short mobility program with foreign universities and research
 431 centers of the Università degli Studi di Napoli Federico II (DR 2243).

432 **Compliance with Ethical Standards**

433 **Conflict of interest** On behalf of all authors, the corresponding author
 434 states that there is no conflict of interest.

435 **Dissimilarity Measure**

436 A general measure $d(\mathbf{x}, \mathbf{y})$ is a dissimilarity measure if the fol-
 437 lowing conditions are verified [33, p.404]:

- 438 1. $d(\mathbf{x}, \mathbf{y}) \geq 0$
- 439 2. $d(\mathbf{x}, \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} = \mathbf{y}$
- 440 3. $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$.

Author Proof

UNCORRECTED PROOF

441 Let $f(x_i; \mu_k, \theta_k)$ be the generic symmetric unimodal multivari-
 442 ate density function of the random variable \mathbf{X} with parameter
 443 θ_k and location parameter μ_k then

444
$$d(x_i, \mu_k) = \log \left(\frac{M_k}{f(x_i; \mu_k, \theta_k)} \right), \tag{23}$$

445 satisfies all the three properties and it is a dissimilarity meas-
 446 ure for $k = 1, \dots, K$.

447 1. $d(x_i, \mu_k) > 0, \forall x_i$.

448 Proof

449
$$0 < \frac{f(x_i; \mu_k, \theta_k)}{M_k} \leq 1 \Rightarrow \frac{M_k}{f(x_i; \mu_k, \theta_k)} \geq 1 \Rightarrow$$

$$\Rightarrow \log \left(\frac{M_k}{f(x_i; \mu_k, \theta_k)} \right) \geq 0.$$

451 2. $d(x_i, \mu_k) = 0 \Leftrightarrow x_i = \mu_k$.

452 2a. $x_i = \mu_k \Rightarrow d(x_i, \mu_k) = 0 \forall x_i$. Proof

453
$$x_i = \mu_k \Rightarrow f(x_i; \mu_k, \theta_k) = f(\mu_k; \mu_k, \theta_k) = M_k \Rightarrow$$

$$\Rightarrow \frac{M_k}{M_k} = 1 \Rightarrow \log(1) = 0,$$

455 2b. $d(x_i, \mu_k) = 0 \Rightarrow x_i = \mu_k, \forall x_i$.

456 Proof

457
$$\log \left(\frac{M_k}{f(x_i; \mu_k, \theta_k)} \right) = 0 \Rightarrow \frac{M_k}{f(x_i; \mu_k, \theta_k)} = 1 \Rightarrow$$

$$\Rightarrow f(x_i; \mu_k, \theta_k) = M_k$$

$$= f(\mu_k; \mu_k, \theta_k) \Rightarrow x_i = \mu_k.$$

459 3. $d(x_i, \mu_k) = d(\mu_k, x_i), \forall x_i$ Proof Given θ_k

460
$$f(x_i; \mu_k, \theta_k) = f(\mu_k; x_i, \theta_k) \Rightarrow \log \left(\frac{M_k}{f(x_i; \mu_k, \theta_k)} \right) = \log \left(\frac{M_k}{f(\mu_k; x_i, \theta_k)} \right)$$

461 □

464 **Addition Details for Data Analyses**

465 **Variable Selection**

466 On each dataset we selected one variable per group using
 467 **1.2.1** hierarchical clustering (Fig. 4).

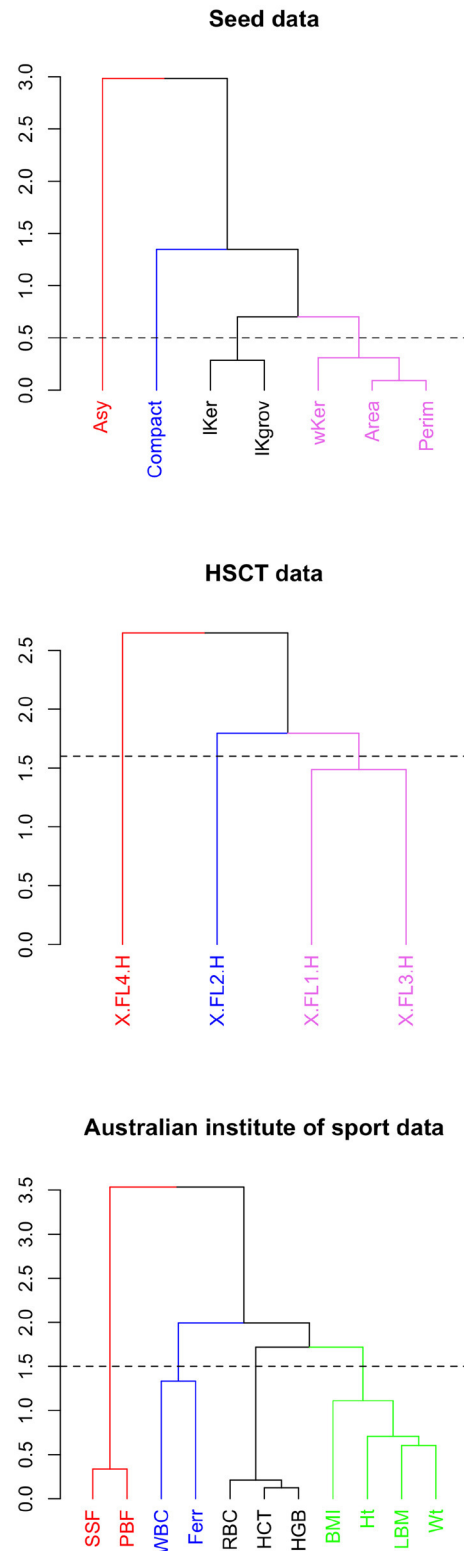


Fig. 4 Variable selection

Author Proof

UNCORRECTED

Simulated Datasets

Table 5 Simulated datasets using multivariate Gaussian distributions Scenarios 1-4

	True		GPDC			TPDC			
μ_1	0.00	0.00	0.00	- 0.02	0.03	- 0.01	0.01	-0.00	-0.01
μ_2	- 7.00	7.00	0.00	- 7.00	7.00	0.01	-7.01	7.01	0.01
μ_3	-7.00	0.00	7.00	-6.98	0.05	6.97	-6.99	0.02	7.02
Σ_1	1.00	0.00	0.00	1.29	-0.29	-0.02	0.97	-0.06	0.01
	0.00	1.00	0.00	-0.29	1.23	0.00	-0.06	0.95	0.00
	0.00	0.00	1.00	-0.02	0.00	0.93	0.01	0.00	0.84
Σ_2	1.00	0.00	0.00	1.04	-0.11	0.00	0.90	-0.02	-0.00
	0.00	1.00	0.00	-0.11	1.09	-0.04	-0.02	0.93	-0.01
	0.00	0.00	1.00	0.00	-0.04	0.90	-0.00	-0.01	0.84
Σ_3	1.00	0.00	0.00	1.04	-0.02	-0.11	0.87	-0.01	-0.01
	0.00	1.00	0.00	-0.02	1.18	-0.28	-0.01	0.89	-0.06
	0.00	0.00	1.00	-0.11	-0.28	1.35	-0.01	-0.06	0.95
μ_1	0.00	0.00	0.00	-0.02	-0.00	0.03	0.02	-0.00	-0.01
μ_2	-7.00	7.00	0.00	-6.99	6.94	0.07	-7.00	6.99	0.04
μ_3	-7.00	0.00	7.00	-6.99	0.01	7.00	-7.00	0.01	7.01
Σ_1	1.00	0.00	0.00	1.28	-0.03	-0.27	0.97	-0.00	-0.04
	0.00	1.00	0.00	-0.03	0.91	0.00	-0.00	0.82	0.00
	0.00	0.00	1.00	-0.27	0.00	1.25	-0.04	0.00	0.97
Σ_2	1.00	0.00	0.00	1.03	-0.12	-0.03	0.85	-0.02	-0.02
	0.00	1.00	0.00	-0.12	1.36	-0.30	-0.02	0.95	-0.07
	0.00	0.00	1.00	-0.03	-0.30	1.23	-0.02	-0.07	0.92
Σ_3	1.00	0.00	0.00	1.03	0.00	-0.12	0.90	0.00	-0.02
	0.00	1.00	0.00	0.00	0.91	-0.03	0.00	0.84	-0.01
	0.00	0.00	1.00	-0.12	-0.03	1.07	-0.02	-0.01	0.90
μ_1	0.00	0.00	0.00	-0.08	0.12	-0.02	0.01	0.03	-0.03
μ_2	-4.00	4.00	0.00	-3.99	4.04	-0.01	-4.02	4.05	-0.01
μ_3	-4.00	0.00	4.00	-3.92	0.22	3.79	-3.95	0.10	3.97
Σ_1	1.00	0.00	0.00	1.51	-0.52	-0.05	1.00	-0.20	-0.01
	0.00	1.00	0.00	-0.52	1.38	-0.00	-0.20	0.94	0.00
	0.00	0.00	1.00	-0.05	-0.00	0.87	-0.01	0.00	0.68
Σ_2	1.00	0.00	0.00	1.01	-0.19	-0.00	0.80	-0.07	-0.00
	0.00	1.00	0.00	-0.19	1.13	-0.05	-0.07	0.87	-0.02
	0.00	0.00	1.00	-0.00	-0.05	0.81	-0.00	-0.02	0.69
Σ_3	1.00	0.00	0.00	1.16	-0.05	-0.26	0.71	-0.02	-0.09
	0.00	1.00	0.00	-0.05	1.50	-0.70	-0.02	0.85	-0.27
	0.00	0.00	1.00	-0.26	-0.70	1.90	-0.09	-0.27	0.99
μ_1	0.00	0.00	0.00	-0.08	-0.01	0.11	0.01	-0.01	0.02
μ_2	-4.00	4.00	0.00	-3.93	3.77	0.24	-3.97	3.94	0.13
μ_3	-4.00	0.00	4.00	-3.98	-0.01	4.03	-4.01	-0.01	4.05
Σ_1	1.00	0.00	0.00	1.51	-0.06	-0.51	0.99	-0.02	-0.19
	0.00	1.00	0.00	-0.06	0.86	-0.01	-0.02	0.67	-0.01
	0.00	0.00	1.00	-0.51	-0.01	1.40	-0.19	-0.01	0.94
Σ_2	1.00	0.00	0.00	1.14	-0.28	-0.05	0.70	-0.11	-0.02
	0.00	1.00	0.00	-0.28	1.91	-0.72	-0.11	1.01	-0.29
	0.00	0.00	1.00	-0.05	-0.72	1.54	-0.02	-0.29	0.88
Σ_3	1.00	0.00	0.00	1.02	-0.00	-0.20	0.78	-0.00	-0.07
	0.00	1.00	0.00	-0.00	0.82	-0.05	-0.00	0.67	-0.02
	0.00	0.00	1.00	-0.20	-0.05	1.12	-0.07	-0.02	0.82

Author Proof

UNCORRECTED PROOF

Table 6 Simulated datasets using multivariate Gaussian distributions Scenarios 5-8

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.20	0.23	-0.03	-0.02	0.04	-0.02
μ_2	-7.00	7.00	0.00	-7.33	8.08	-0.46	-7.28	7.85	-0.40
μ_3	-7.00	0.00	7.00	-7.09	0.86	6.33	-7.03	0.40	6.78
Σ_1	1.00	0.00	0.00	2.68	-0.93	0.09	0.90	-0.09	0.02
	0.00	1.00	0.00	-0.93	2.60	-0.06	-0.09	0.89	-0.00
	0.00	0.00	1.00	0.09	-0.06	1.69	0.02	-0.00	0.74
Σ_2	16.00	0.00	0.00	12.84	0.95	0.09	11.36	0.46	-0.00
	0.00	16.00	0.00	0.95	13.43	1.24	0.46	12.01	0.80
	0.00	0.00	16.00	0.09	1.24	11.82	-0.00	0.80	10.36
Σ_3	2.25	0.00	0.00	4.40	-0.23	0.39	2.14	-0.09	0.13
	0.00	2.25	0.00	-0.23	6.37	-2.79	-0.09	2.79	-0.85
	0.00	0.00	2.25	0.39	-2.79	7.28	0.13	-0.85	2.94
μ_1	0.00	0.00	0.00	-0.12	0.06	0.01	-0.00	0.01	-0.02
μ_2	-7.00	7.00	0.00	-7.10	7.23	0.80	-7.07	7.11	0.87
μ_3	-7.00	0.00	7.00	-7.02	0.06	7.00	-7.01	0.01	7.06
Σ_1	1.00	0.00	0.00	2.09	-0.29	-0.27	0.91	-0.04	-0.02
	0.00	1.00	0.00	-0.29	1.31	-0.11	-0.04	0.73	-0.02
	0.00	0.00	1.00	-0.27	-0.11	1.76	-0.02	-0.02	0.87
Σ_2	16.00	0.00	0.00	11.06	-0.17	-0.08	9.14	-0.44	-0.06
	0.00	16.00	0.00	-0.17	17.36	-4.90	-0.44	14.91	-4.73
	0.00	0.00	16.00	-0.08	-4.90	14.27	-0.06	-4.73	11.95
Σ_3	2.25	0.00	0.00	2.42	-0.03	-0.02	1.86	-0.01	-0.02
	0.00	2.25	0.00	-0.03	2.29	-0.44	-0.01	1.70	-0.09
	0.00	0.00	2.25	-0.02	-0.44	3.02	-0.02	-0.09	2.01
μ_1	0.00	0.00	0.00	-0.07	0.17	-0.03	0.02	0.04	-0.03
μ_2	-4.00	4.00	0.00	-4.65	5.39	-0.79	-4.55	5.17	-0.65
μ_3	-4.00	0.00	4.00	-4.12	0.66	3.84	-4.07	0.39	4.00
Σ_1	1.00	0.00	0.00	2.67	-0.07	0.01	1.04	0.01	0.00
	0.00	1.00	0.00	-0.07	2.53	-0.04	0.01	1.02	-0.01
	0.00	0.00	1.00	0.01	-0.04	2.18	0.00	-0.01	0.93
Σ_2	16.00	0.00	0.00	12.16	1.15	-0.07	9.85	0.67	-0.06
	0.00	16.00	0.00	1.15	13.16	1.16	0.67	11.00	0.70
	0.00	0.00	16.00	-0.07	1.16	11.51	-0.06	0.70	9.30
Σ_3	2.25	0.00	0.00	5.01	0.07	0.18	2.59	-0.01	0.04
	0.00	2.25	0.00	0.07	5.12	-0.24	-0.01	2.60	-0.05
	0.00	0.00	2.25	0.18	-0.24	6.08	0.04	-0.05	2.92
μ_1	0.00	0.00	0.00	-0.31	0.16	-0.01	-0.18	0.10	0.03
μ_2	-4.00	4.00	0.00	-3.84	2.77	2.35	-1.30	0.33	1.09
μ_3	-4.00	0.00	4.00	-4.18	-0.23	3.82	-4.21	-0.25	3.86
Σ_1	1.00	0.00	0.00	3.00	-0.40	-0.15	1.75	-0.20	-0.20
	0.00	1.00	0.00	-0.40	2.29	-0.32	-0.20	1.40	-0.11
	0.00	0.00	1.00	-0.15	-0.32	2.80	-0.20	-0.11	1.61
Σ_2	16.00	0.00	0.00	5.58	-0.02	0.27	1.99	-0.17	-0.12
	0.00	16.00	0.00	-0.02	9.10	-3.30	-0.17	1.88	-0.38
	0.00	0.00	16.00	0.27	-3.30	6.61	-0.12	-0.38	2.05
Σ_3	2.25	0.00	0.00	2.72	-0.01	0.18	2.01	-0.05	0.13
	0.00	2.25	0.00	-0.01	2.83	-0.61	-0.05	2.04	-0.47
	0.00	0.00	2.25	0.18	-0.61	3.37	0.13	-0.47	2.34

Author Proof

UNCORRECTED PROOF

Table 7 Simulated datasets using multivariate Gaussian distributions Scenarios 9-12

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.04	0.04	0.02	0.01	0.00	0.02
μ_2	-7.00	7.00	0.00	-7.01	7.00	0.03	-7.02	7.01	0.02
μ_3	-7.00	0.00	7.00	-6.98	0.03	7.02	-6.97	0.03	7.03
Σ_1	1.00	0.00	0.00	1.36	-0.34	-0.01	1.00	-0.07	0.00
	0.00	1.00	0.00	-0.34	1.30	-0.01	-0.07	0.98	-0.00
	0.00	0.00	1.00	-0.01	-0.01	0.96	0.00	-0.00	0.86
Σ_2	1.00	-0.50	-0.50	1.07	-0.60	-0.47	0.92	-0.47	-0.43
	-0.50	1.00	0.50	-0.60	1.11	0.46	-0.47	0.94	0.43
	-0.50	0.50	1.00	-0.47	0.46	0.91	-0.43	0.43	0.85
Σ_3	1.00	0.70	0.70	0.93	0.63	0.64	0.86	0.59	0.61
	0.70	1.00	0.70	0.63	0.93	0.62	0.59	0.85	0.60
	0.70	0.70	1.00	0.64	0.62	0.98	0.61	0.60	0.89
μ_1	0.00	0.00	0.00	-0.04	-0.00	0.05	0.00	-0.00	0.01
μ_2	-7.00	7.00	0.00	-7.00	6.97	0.05	-7.02	7.00	0.04
μ_3	-7.00	0.00	7.00	-6.99	0.02	7.01	-6.99	0.02	7.01
Σ_1	1.00	0.00	0.00	1.40	-0.01	-0.35	1.01	-0.00	-0.06
	0.00	1.00	0.00	-0.01	0.97	0.01	-0.00	0.88	0.00
	0.00	0.00	1.00	-0.35	0.01	1.36	-0.06	0.00	1.01
Σ_2	1.00	-0.50	-0.50	1.06	-0.59	-0.50	0.90	-0.46	-0.46
	-0.50	1.00	0.50	-0.59	1.19	0.39	-0.46	0.93	0.43
	-0.50	0.50	1.00	-0.50	0.39	1.06	-0.46	0.43	0.91
Σ_3	1.00	0.70	0.70	0.94	0.63	0.64	0.88	0.60	0.61
	0.70	1.00	0.70	0.63	0.88	0.63	0.60	0.84	0.60
	0.70	0.70	1.00	0.64	0.63	0.94	0.61	0.60	0.88
μ_1	0.00	0.00	0.00	-0.19	0.22	0.04	-0.09	0.11	0.05
μ_2	-4.00	4.00	0.00	-4.05	4.04	0.08	-4.06	4.06	0.07
μ_3	-4.00	0.00	4.00	-3.95	0.05	4.04	-3.94	0.05	4.08
Σ_1	1.00	0.00	0.00	1.78	-0.76	-0.00	1.32	-0.41	-0.01
	0.00	1.00	0.00	-0.76	1.62	-0.08	-0.41	1.20	-0.05
	0.00	0.00	1.00	-0.00	-0.08	0.97	-0.01	-0.05	0.83
Σ_2	1.00	-0.50	-0.50	1.05	-0.64	-0.44	0.84	-0.47	-0.38
	-0.50	1.00	0.50	-0.64	1.17	0.43	-0.47	0.91	0.37
	-0.50	0.50	1.00	-0.44	0.43	0.83	-0.38	0.37	0.72
Σ_3	1.00	0.70	0.70	0.86	0.57	0.57	0.57	0.40	0.42
	0.70	1.00	0.70	0.57	0.90	0.55	0.40	0.59	0.42
	0.70	0.70	1.00	0.57	0.55	0.99	0.42	0.42	0.65
μ_1	0.00	0.00	0.00	-0.17	-0.00	0.24	-0.08	0.01	0.12
μ_2	-4.00	4.00	0.00	-3.94	3.85	0.12	-4.03	4.02	0.08
μ_3	-4.00	0.00	4.00	-3.99	0.03	4.05	-3.99	0.03	4.04
Σ_1	1.00	0.00	0.00	1.89	0.04	-0.79	1.34	-0.02	-0.40
	0.00	1.00	0.00	0.04	0.97	0.02	-0.02	0.85	0.00
	0.00	0.00	1.00	-0.79	0.02	1.70	-0.40	0.00	1.25
Σ_2	1.00	-0.50	-0.50	1.19	-0.79	-0.46	0.68	-0.40	-0.32
	-0.50	1.00	0.50	-0.79	1.66	0.22	-0.40	0.81	0.28
	-0.50	0.50	1.00	-0.46	0.22	1.23	-0.32	0.28	0.70
Σ_3	1.00	0.70	0.70	0.88	0.59	0.62	0.78	0.53	0.56
	0.70	1.00	0.70	0.59	0.77	0.59	0.53	0.70	0.53
	0.70	0.70	1.00	0.62	0.59	0.91	0.56	0.53	0.81

Author Proof

Table 8 Simulated datasets using multivariate Gaussian distributions Scenarios 13–16

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.44	0.48	-0.24	-0.09	0.11	-0.06
μ_2	-7.00	7.00	0.00	-7.94	8.25	0.64	-7.73	7.75	0.58
μ_3	-7.00	0.00	7.00	-6.98	0.63	6.52	-6.92	0.09	7.09
Σ_1	1.00	0.00	0.00	4.11	-2.42	0.32	1.24	-0.36	0.10
	0.00	1.00	0.00	-2.42	3.92	-0.39	-0.36	1.20	-0.11
	0.00	0.00	1.00	0.32	-0.39	2.86	0.10	-0.11	1.06
Σ_2	16.00	-8.00	-8.00	13.01	-6.32	-6.60	12.04	-5.73	-5.83
	-8.00	16.00	8.00	-6.32	13.59	6.70	-5.73	12.79	5.60
	-8.00	8.00	16.00	-6.60	6.70	11.52	-5.83	5.60	11.29
Σ_3	2.25	1.57	1.57	3.11	0.73	0.83	1.29	0.88	0.93
	1.57	2.25	1.57	0.73	5.46	-0.75	0.88	1.37	0.90
	1.57	1.57	2.25	0.83	-0.75	6.25	0.93	0.90	1.46
μ_1	0.00	0.00	0.00	-0.16	0.10	0.00	-0.01	0.02	-0.01
μ_2	-7.00	7.00	0.00	-7.72	7.58	0.98	-7.61	7.51	0.92
μ_3	-7.00	0.00	7.00	-6.96	0.05	7.05	-6.96	0.05	7.06
Σ_1	1.00	0.00	0.00	2.40	-0.40	-0.55	0.94	-0.05	-0.04
	0.00	1.00	0.00	-0.40	1.46	-0.08	-0.05	0.76	-0.03
	0.00	0.00	1.00	-0.55	-0.08	2.40	-0.04	-0.03	0.96
Σ_2	16.00	-8.00	-8.00	12.21	-5.92	-5.73	10.56	-5.29	-4.72
	-8.00	16.00	8.00	-5.92	15.67	3.66	-5.29	13.38	2.96
	-8.00	8.00	16.00	-5.73	3.66	13.43	-4.72	2.96	11.09
Σ_3	2.25	1.57	1.57	2.00	1.32	1.37	1.80	1.21	1.27
	1.57	2.25	1.57	1.32	1.87	1.29	1.21	1.66	1.20
	1.57	1.57	2.25	1.37	1.29	2.11	1.27	1.20	1.84
μ_1	0.00	0.00	0.00	-0.25	0.41	-0.21	-0.21	0.17	0.02
μ_2	-4.00	4.00	0.00	-4.99	5.11	0.35	-5.06	4.94	0.45
μ_3	-4.00	0.00	4.00	-4.12	1.09	3.35	-3.52	0.39	3.62
Σ_1	1.00	0.00	0.00	3.91	-1.63	-1.06	2.06	-0.51	-0.61
	0.00	1.00	0.00	-1.63	3.78	0.90	-0.51	1.79	0.34
	0.00	0.00	1.00	-1.06	0.90	3.57	-0.61	0.34	2.05
Σ_2	16.00	-8.00	-8.00	11.09	-5.47	-5.53	10.45	-5.07	-5.30
	-8.00	16.00	8.00	-5.47	12.68	6.01	-5.07	11.75	5.11
	-8.00	8.00	16.00	-5.53	6.01	10.45	-5.30	5.11	10.05
Σ_3	2.25	1.57	1.57	5.54	-1.45	-0.85	1.76	0.22	0.27
	1.57	2.25	1.57	-1.45	6.76	1.81	0.22	1.96	0.84
	1.57	1.57	2.25	-0.85	1.81	6.94	0.27	0.84	2.03
μ_1	0.00	0.00	0.00	-0.37	0.29	-0.14	-0.32	0.32	-0.02
μ_2	-4.00	4.00	0.00	-2.03	1.59	1.47	-0.32	0.32	-0.02
μ_3	-4.00	0.00	4.00	-4.66	-0.55	3.40	-4.41	-0.34	3.66
Σ_1	1.00	0.00	0.00	3.80	-1.26	-1.06	2.70	-1.22	-0.64
	0.00	1.00	0.00	-1.26	2.89	0.49	-1.22	2.39	0.36
	0.00	0.00	1.00	-1.06	0.49	3.82	-0.64	0.36	2.35
Σ_2	16.00	-8.00	-8.00	4.97	-2.68	-0.61	2.70	-1.22	-0.64
	-8.00	16.00	8.00	-2.68	6.30	-0.13	-1.22	2.39	0.36
	-8.00	8.00	16.00	-0.61	-0.13	5.50	-0.64	0.36	2.35
Σ_3	2.25	1.57	1.57	1.92	0.60	0.88	1.20	0.45	0.62
	1.57	2.25	1.57	0.60	2.02	0.78	0.45	1.18	0.50
	1.57	1.57	2.25	0.88	0.78	2.02	0.62	0.50	1.27

Author Proof

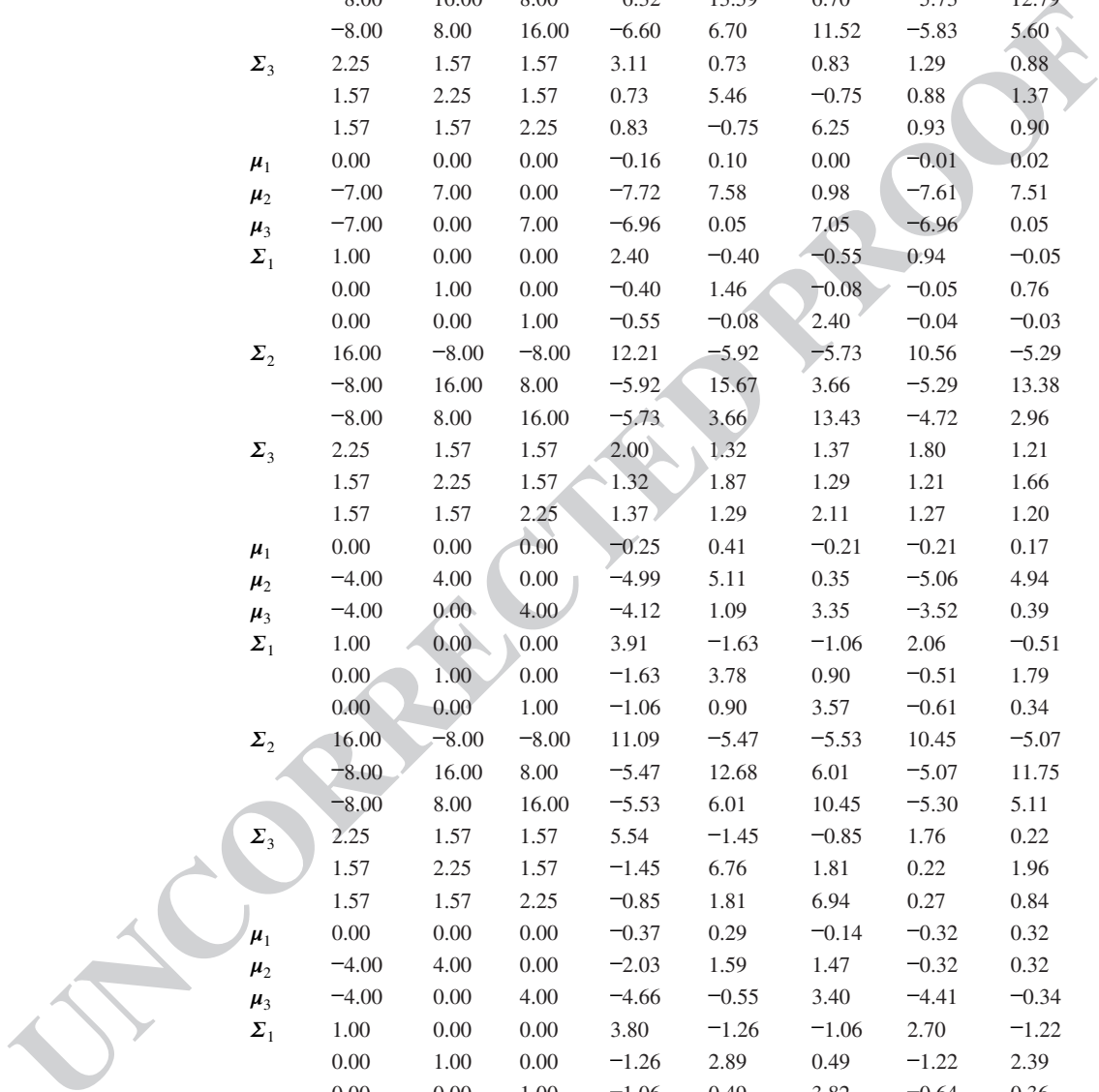


Table 9 Simulated datasets using multivariate Student-*t* distributions Scenarios 1–4

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.05	0.06	-0.01	0.01	-0.00	-0.01
μ_2	-7.00	7.00	0.00	-6.98	7.02	-0.00	-7.00	7.02	0.00
μ_3	-7.00	0.00	7.00	-6.96	0.11	6.90	-6.99	0.01	7.02
Σ_1	1.00	0.00	0.00	2.02	-0.60	-0.08	0.96	-0.09	-0.03
	0.00	1.00	0.00	-0.60	1.94	-0.01	-0.09	0.94	-0.01
	0.00	0.00	1.00	-0.08	-0.01	1.27	-0.03	-0.01	0.76
Σ_2	1.00	0.00	0.00	1.47	-0.22	-0.01	1.01	-0.04	-0.01
	0.00	1.00	0.00	-0.22	1.62	-0.04	-0.04	1.06	-0.00
	0.00	0.00	1.00	-0.01	-0.04	1.24	-0.01	-0.00	0.90
Σ_3	1.00	0.00	0.00	1.65	-0.06	-0.32	0.81	-0.03	-0.04
	0.00	1.00	0.00	-0.06	2.13	-0.72	-0.03	0.88	-0.08
	0.00	0.00	1.00	-0.32	-0.72	2.59	-0.04	-0.08	0.97
μ_1	0.00	0.00	0.00	-0.06	-0.01	0.07	0.01	-0.01	-0.00
μ_2	-7.00	7.00	0.00	-6.96	6.87	0.11	-7.02	7.00	0.01
μ_3	-7.00	0.00	7.00	-7.00	-0.02	7.02	-7.02	-0.01	7.03
Σ_1	1.00	0.00	0.00	2.05	-0.06	-0.63	0.97	-0.02	-0.11
	0.00	1.00	0.00	-0.06	1.28	-0.01	-0.02	0.78	-0.00
	0.00	0.00	1.00	-0.63	-0.01	1.92	-0.11	-0.00	0.95
Σ_2	1.00	0.00	0.00	1.67	-0.34	-0.01	0.80	-0.05	0.00
	0.00	1.00	0.00	-0.34	2.51	-0.73	-0.05	0.95	-0.10
	0.00	0.00	1.00	-0.01	-0.73	2.07	0.00	-0.10	0.90
Σ_3	1.00	0.00	0.00	1.52	-0.04	-0.25	1.03	-0.03	-0.06
	0.00	1.00	0.00	-0.04	1.24	-0.05	-0.03	0.92	-0.02
	0.00	0.00	1.00	-0.25	-0.05	1.60	-0.06	-0.02	1.05
μ_1	0.00	0.00	0.00	-0.09	0.13	-0.03	-0.01	0.05	-0.02
μ_2	-4.00	4.00	0.00	-3.99	4.06	-0.03	-4.02	4.07	-0.02
μ_3	-4.00	0.00	4.00	-3.88	0.29	3.71	-3.94	0.14	3.89
Σ_1	1.00	0.00	0.00	1.98	-0.61	-0.08	0.95	-0.21	-0.03
	0.00	1.00	0.00	-0.61	1.86	-0.01	-0.21	0.91	-0.01
	0.00	0.00	1.00	-0.08	-0.01	1.19	-0.03	-0.01	0.64
Σ_2	1.00	0.00	0.00	1.33	-0.20	-0.01	0.75	-0.05	-0.01
	0.00	1.00	0.00	-0.20	1.54	-0.03	-0.05	0.81	0.00
	0.00	0.00	1.00	-0.01	-0.03	1.11	-0.01	0.00	0.64
Σ_3	1.00	0.00	0.00	1.70	-0.09	-0.40	0.86	-0.05	-0.16
	0.00	1.00	0.00	-0.09	2.22	-0.88	-0.05	1.09	-0.37
	0.00	0.00	1.00	-0.40	-0.88	2.85	-0.16	-0.37	1.32
μ_1	0.00	0.00	0.00	-0.10	-0.03	0.13	-0.01	-0.02	0.04
μ_2	-4.00	4.00	0.00	-3.90	3.67	0.28	-3.97	3.87	0.14
μ_3	-4.00	0.00	4.00	-4.00	-0.05	4.06	-4.03	-0.03	4.08
Σ_1	1.00	0.00	0.00	2.01	-0.05	-0.64	0.96	-0.02	-0.23
	0.00	1.00	0.00	-0.05	1.20	-0.02	-0.02	0.64	-0.01
	0.00	0.00	1.00	-0.64	-0.02	1.82	-0.23	-0.01	0.91
Σ_2	1.00	0.00	0.00	1.72	-0.40	-0.06	0.83	-0.15	-0.02
	0.00	1.00	0.00	-0.40	2.75	-0.90	-0.15	1.25	-0.39
	0.00	0.00	1.00	-0.06	-0.90	2.18	-0.02	-0.39	1.07
Σ_3	1.00	0.00	0.00	1.38	-0.03	-0.23	0.78	-0.02	-0.08
	0.00	1.00	0.00	-0.03	1.12	-0.03	-0.02	0.66	-0.01
	0.00	0.00	1.00	-0.23	-0.03	1.53	-0.08	-0.01	0.82

Table 10 Simulated datasets using multivariate Student-*t* distributions Scenarios 5–8

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.16	0.17	-0.02	-0.01	0.02	-0.01
μ_2	-7.00	7.00	0.00	-7.04	7.27	-0.07	-7.04	7.18	-0.05
μ_3	-7.00	0.00	7.00	-6.99	0.47	6.59	-7.00	0.14	6.93
Σ_1	1.00	0.00	0.00	2.80	-1.12	-0.04	0.97	-0.15	-0.02
	0.00	1.00	0.00	-1.12	2.72	-0.05	-0.15	0.96	-0.01
	0.00	0.00	1.00	-0.04	-0.05	1.60	-0.02	-0.01	0.74
Σ_2	16.00	0.00	0.00	4.83	-0.20	0.02	3.58	-0.14	-0.02
	0.00	16.00	0.00	-0.20	5.39	0.09	-0.14	3.84	0.06
	0.00	0.00	16.00	0.02	0.09	4.25	-0.02	0.06	3.14
Σ_3	2.25	0.00	0.00	2.90	-0.14	-0.21	1.32	-0.07	-0.05
	0.00	2.25	0.00	-0.14	4.69	-2.42	-0.07	1.77	-0.56
	0.00	0.00	2.25	-0.21	-2.42	5.80	-0.05	-0.56	1.97
μ_1	0.00	0.00	0.00	-0.10	0.00	0.08	0.00	-0.01	-0.00
μ_2	-7.00	7.00	0.00	-6.96	6.76	0.36	-7.01	6.93	0.18
μ_3	-7.00	0.00	7.00	-7.02	-0.02	7.04	-7.03	-0.02	7.06
Σ_1	1.00	0.00	0.00	2.32	-0.18	-0.65	0.91	-0.03	-0.10
	0.00	1.00	0.00	-0.18	1.42	-0.03	-0.03	0.71	-0.01
	0.00	0.00	1.00	-0.65	-0.03	2.06	-0.10	-0.01	0.87
Σ_2	16.00	0.00	0.00	5.13	-0.65	-0.00	2.94	-0.34	-0.02
	0.00	16.00	0.00	-0.65	8.04	-2.46	-0.34	4.20	-1.14
	0.00	0.00	16.00	-0.00	-2.46	6.76	-0.02	-1.14	3.78
Σ_3	2.25	0.00	0.00	2.13	-0.05	-0.22	1.48	-0.04	-0.08
	0.00	2.25	0.00	-0.05	1.83	-0.20	-0.04	1.32	-0.04
	0.00	0.00	2.25	-0.22	-0.20	2.40	-0.08	-0.04	1.56
μ_1	0.00	0.00	0.00	-0.16	0.21	-0.05	-0.05	0.09	-0.03
μ_2	-4.00	4.00	0.00	-4.12	4.46	-0.19	-4.12	4.38	-0.15
μ_3	-4.00	0.00	4.00	-3.99	0.62	3.53	-3.98	0.40	3.72
Σ_1	1.00	0.00	0.00	2.57	-0.62	-0.02	1.11	-0.20	-0.02
	0.00	1.00	0.00	-0.62	2.48	-0.07	-0.20	1.08	-0.02
	0.00	0.00	1.00	-0.02	-0.07	1.80	-0.02	-0.02	0.82
Σ_2	16.00	0.00	0.00	4.42	-0.02	0.03	2.61	-0.05	-0.01
	0.00	16.00	0.00	-0.02	5.15	0.19	-0.05	2.91	0.10
	0.00	0.00	16.00	0.03	0.19	4.06	-0.01	0.10	2.34
Σ_3	2.25	0.00	0.00	3.18	-0.14	-0.17	1.56	-0.10	-0.10
	0.00	2.25	0.00	-0.14	3.89	-1.26	-0.10	1.96	-0.63
	0.00	0.00	2.25	-0.17	-1.26	4.89	-0.10	-0.63	2.30
μ_1	0.00	0.00	0.00	-0.16	0.01	0.10	-0.05	-0.00	0.03
μ_2	-4.00	4.00	0.00	-3.95	3.53	0.75	-3.95	3.45	0.75
μ_3	-4.00	0.00	4.00	-4.06	-0.05	4.11	-4.06	-0.05	4.14
Σ_1	1.00	0.00	0.00	2.40	-0.17	-0.49	1.08	-0.05	-0.21
	0.00	1.00	0.00	-0.17	1.57	-0.12	-0.05	0.76	-0.04
	0.00	0.00	1.00	-0.49	-0.12	2.15	-0.21	-0.04	1.02
Σ_2	16.00	0.00	0.00	4.25	-0.42	-0.04	2.83	-0.38	-0.09
	0.00	16.00	0.00	-0.42	6.72	-2.02	-0.38	4.79	-1.79
	0.00	0.00	16.00	-0.04	-2.02	5.15	-0.09	-1.79	3.74
Σ_3	2.25	0.00	0.00	2.04	-0.05	-0.13	1.10	-0.03	-0.07
	0.00	2.25	0.00	-0.05	1.81	-0.16	-0.03	0.97	-0.02
	0.00	0.00	2.25	-0.13	-0.16	2.41	-0.07	-0.02	1.21

Author Proof

UNCORRECTED PROOF

Table 11 Simulated datasets using multivariate Student-*t* distributions Scenarios 9–12

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.11	0.11	0.04	-0.01	0.01	0.03
μ_2	-7.00	7.00	0.00	-7.01	7.02	0.04	-7.02	7.03	0.03
μ_3	-7.00	0.00	7.00	-6.97	0.03	7.03	-6.98	0.02	7.03
Σ_1	1.00	0.00	0.00	2.37	-0.83	-0.08	1.16	-0.14	-0.04
	0.00	1.00	0.00	-0.83	2.26	-0.06	-0.14	1.12	-0.02
	0.00	0.00	1.00	-0.08	-0.06	1.43	-0.04	-0.02	0.91
Σ_2	1.00	-0.50	-0.50	1.55	-0.92	-0.69	1.07	-0.57	-0.51
	-0.50	1.00	0.50	-0.92	1.69	0.68	-0.57	1.12	0.51
	-0.50	0.50	1.00	-0.69	0.68	1.31	-0.51	0.51	0.97
Σ_3	1.00	0.70	0.70	1.32	0.89	0.89	0.77	0.53	0.56
	0.70	1.00	0.70	0.89	1.44	0.89	0.53	0.80	0.57
	0.70	0.70	1.00	0.89	0.89	1.54	0.56	0.57	0.87
μ_1	0.00	0.00	0.00	-0.10	-0.01	0.14	-0.01	-0.01	0.02
μ_2	-7.00	7.00	0.00	-6.98	6.92	0.06	-7.03	7.01	0.02
μ_3	-7.00	0.00	7.00	-7.01	-0.00	7.01	-7.01	-0.00	7.01
Σ_1	1.00	0.00	0.00	2.46	-0.00	-0.87	1.11	-0.02	-0.14
	0.00	1.00	0.00	-0.00	1.47	0.03	-0.02	0.91	-0.00
	0.00	0.00	1.00	-0.87	0.03	2.26	-0.14	-0.00	1.07
Σ_2	1.00	-0.50	-0.50	1.71	-1.04	-0.66	0.83	-0.45	-0.40
	-0.50	1.00	0.50	-1.04	2.24	0.43	-0.45	0.92	0.39
	-0.50	0.50	1.00	-0.66	0.43	1.72	-0.40	0.39	0.86
Σ_3	1.00	0.70	0.70	1.32	0.86	0.89	1.04	0.69	0.72
	0.70	1.00	0.70	0.86	1.18	0.86	0.69	0.97	0.70
	0.70	0.70	1.00	0.89	0.86	1.34	0.72	0.70	1.06
μ_1	0.00	0.00	0.00	-0.26	0.29	0.02	-0.13	0.14	0.03
μ_2	-4.00	4.00	0.00	-4.08	4.09	0.09	-4.07	4.08	0.07
μ_3	-4.00	0.00	4.00	-3.93	0.07	4.01	-3.96	0.02	4.06
Σ_1	1.00	0.00	0.00	2.49	-1.03	-0.08	1.27	-0.44	-0.04
	0.00	1.00	0.00	-1.03	2.30	-0.07	-0.44	1.19	-0.06
	0.00	0.00	1.00	-0.08	-0.07	1.47	-0.04	-0.06	0.84
Σ_2	1.00	-0.50	-0.50	1.36	-0.80	-0.62	0.89	-0.49	-0.41
	-0.50	1.00	0.50	-0.80	1.55	0.61	-0.49	0.98	0.42
	-0.50	0.50	1.00	-0.62	0.61	1.14	-0.41	0.42	0.77
Σ_3	1.00	0.70	0.70	1.25	0.77	0.73	0.60	0.42	0.44
	0.70	1.00	0.70	0.77	1.43	0.76	0.42	0.65	0.46
	0.70	0.70	1.00	0.73	0.76	1.62	0.44	0.46	0.73
μ_1	0.00	0.00	0.00	-0.28	-0.04	0.36	-0.11	-0.01	0.18
μ_2	-4.00	4.00	0.00	-3.89	3.71	0.18	-4.01	3.97	0.04
μ_3	-4.00	0.00	4.00	-4.01	-0.00	4.06	-4.01	0.01	4.05
Σ_1	1.00	0.00	0.00	2.81	0.18	-1.16	1.36	0.03	-0.50
	0.00	1.00	0.00	0.18	1.43	0.11	0.03	0.82	0.03
	0.00	0.00	1.00	-1.16	0.11	2.44	-0.50	0.03	1.25
Σ_2	1.00	-0.50	-0.50	1.84	-1.10	-0.55	0.81	-0.49	-0.35
	-0.50	1.00	0.50	-1.10	2.72	0.13	-0.49	1.02	0.30
	-0.50	0.50	1.00	-0.55	0.13	2.07	-0.35	0.30	0.84
Σ_3	1.00	0.70	0.70	1.19	0.79	0.87	0.86	0.57	0.63
	0.70	1.00	0.70	0.79	1.00	0.81	0.57	0.76	0.59
	0.70	0.70	1.00	0.87	0.81	1.25	0.63	0.59	0.91

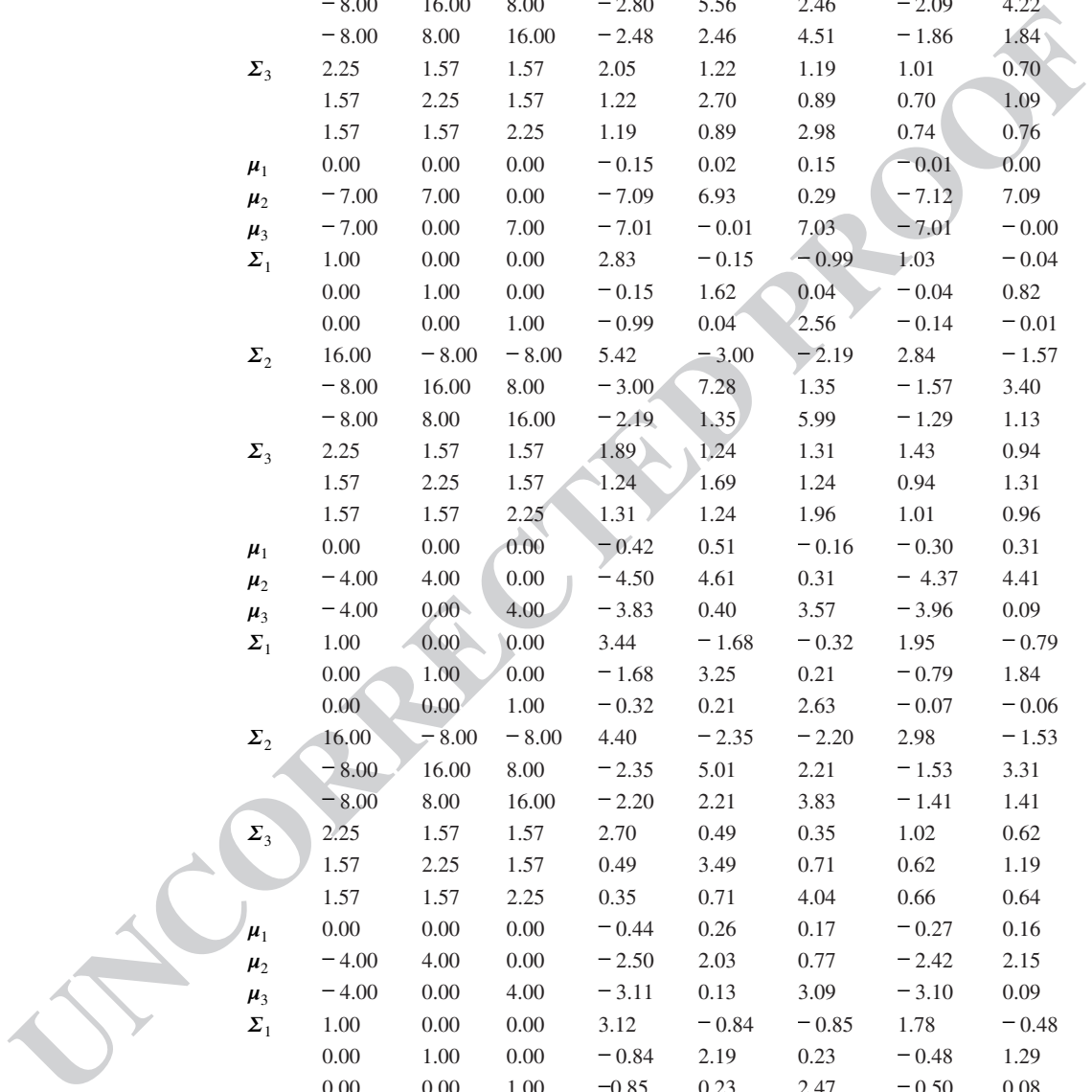
Author Proof

UNCORRECTED PROOF

Table 12 Simulated datasets using multivariate Student-*t* distributions Scenarios 13–16

	True			GPDC			TPDC		
μ_1	0.00	0.00	0.00	-0.35	0.35	-0.04	-0.06	0.06	0.01
μ_2	-7.00	7.00	0.00	-7.22	7.28	0.20	-7.15	7.17	0.15
μ_3	-7.00	0.00	7.00	-6.97	0.12	6.95	-6.97	0.02	7.05
Σ_1	1.00	0.00	0.00	3.87	-2.17	-0.03	1.22	-0.33	-0.00
	0.00	1.00	0.00	-2.17	3.67	-0.13	-0.33	1.18	-0.06
	0.00	0.00	1.00	-0.03	-0.13	2.15	-0.00	-0.06	0.90
Σ_2	16.00	-8.00	-8.00	5.07	-2.80	-2.48	3.94	-2.09	-1.86
	-8.00	16.00	8.00	-2.80	5.56	2.46	-2.09	4.22	1.84
	-8.00	8.00	16.00	-2.48	2.46	4.51	-1.86	1.84	3.55
Σ_3	2.25	1.57	1.57	2.05	1.22	1.19	1.01	0.70	0.74
	1.57	2.25	1.57	1.22	2.70	0.89	0.70	1.09	0.76
	1.57	1.57	2.25	1.19	0.89	2.98	0.74	0.76	1.18
μ_1	0.00	0.00	0.00	-0.15	0.02	0.15	-0.01	0.00	0.02
μ_2	-7.00	7.00	0.00	-7.09	6.93	0.29	-7.12	7.09	0.13
μ_3	-7.00	0.00	7.00	-7.01	-0.01	7.03	-7.01	-0.00	7.02
Σ_1	1.00	0.00	0.00	2.83	-0.15	-0.99	1.03	-0.04	-0.14
	0.00	1.00	0.00	-0.15	1.62	0.04	-0.04	0.82	-0.01
	0.00	0.00	1.00	-0.99	0.04	2.56	-0.14	-0.01	0.99
Σ_2	16.00	-8.00	-8.00	5.42	-3.00	-2.19	2.84	-1.57	-1.29
	-8.00	16.00	8.00	-3.00	7.28	1.35	-1.57	3.40	1.13
	-8.00	8.00	16.00	-2.19	1.35	5.99	-1.29	1.13	3.04
Σ_3	2.25	1.57	1.57	1.89	1.24	1.31	1.43	0.94	1.01
	1.57	2.25	1.57	1.24	1.69	1.24	0.94	1.31	0.96
	1.57	1.57	2.25	1.31	1.24	1.96	1.01	0.96	1.48
μ_1	0.00	0.00	0.00	-0.42	0.51	-0.16	-0.30	0.31	-0.08
μ_2	-4.00	4.00	0.00	-4.50	4.61	0.31	-4.37	4.41	0.34
μ_3	-4.00	0.00	4.00	-3.83	0.40	3.57	-3.96	0.09	4.05
Σ_1	1.00	0.00	0.00	3.44	-1.68	-0.32	1.95	-0.79	-0.07
	0.00	1.00	0.00	-1.68	3.25	0.21	-0.79	1.84	-0.06
	0.00	0.00	1.00	-0.32	0.21	2.63	-0.07	-0.06	1.50
Σ_2	16.00	-8.00	-8.00	4.40	-2.35	-2.20	2.98	-1.53	-1.41
	-8.00	16.00	8.00	-2.35	5.01	2.21	-1.53	3.31	1.41
	-8.00	8.00	16.00	-2.20	2.21	3.83	-1.41	1.41	2.68
Σ_3	2.25	1.57	1.57	2.70	0.49	0.35	1.02	0.62	0.66
	1.57	2.25	1.57	0.49	3.49	0.71	0.62	1.19	0.64
	1.57	1.57	2.25	0.35	0.71	4.04	0.66	0.64	1.35
μ_1	0.00	0.00	0.00	-0.44	0.26	0.17	-0.27	0.16	0.13
μ_2	-4.00	4.00	0.00	-2.50	2.03	0.77	-2.42	2.15	0.36
μ_3	-4.00	0.00	4.00	-3.11	0.13	3.09	-3.10	0.09	3.20
Σ_1	1.00	0.00	0.00	3.12	-0.84	-0.85	1.78	-0.48	-0.50
	0.00	1.00	0.00	-0.84	2.19	0.23	-0.48	1.29	0.08
	0.00	0.00	1.00	-0.85	0.23	2.47	-0.50	0.08	1.41
Σ_2	16.00	-8.00	-8.00	3.66	-2.07	-0.63	2.46	-1.40	-0.64
	-8.00	16.00	8.00	-2.07	4.99	0.16	-1.40	3.27	0.14
	-8.00	8.00	16.00	-0.63	0.16	3.83	-0.64	0.14	2.48
Σ_3	2.25	1.57	1.57	1.49	0.66	0.82	0.90	0.52	0.63
	1.57	2.25	1.57	0.66	1.53	0.57	0.52	0.83	0.50
	1.57	1.57	2.25	0.82	0.57	1.82	0.63	0.50	1.03

Author Proof



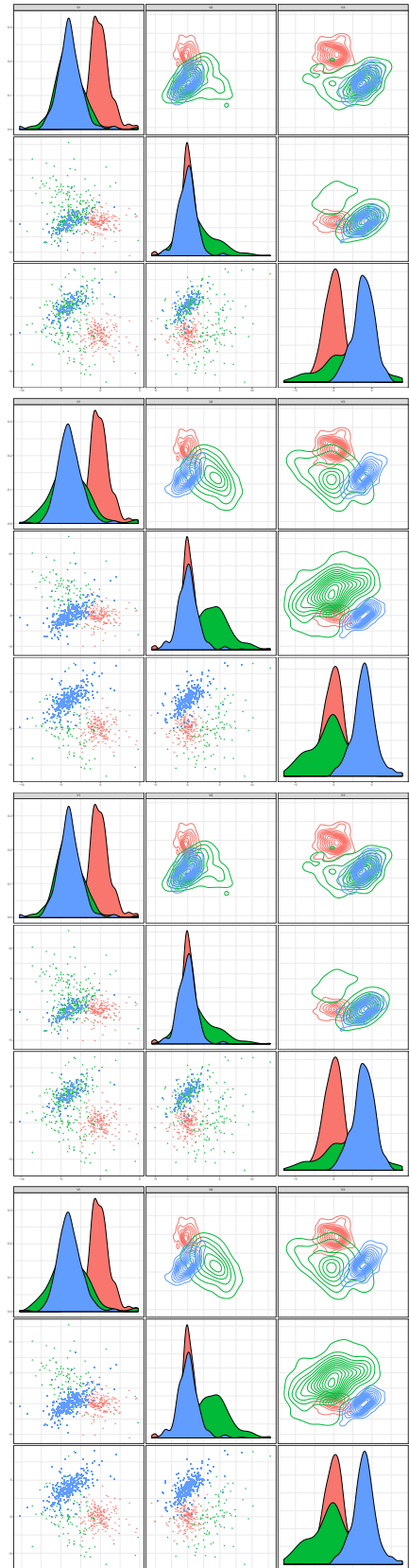
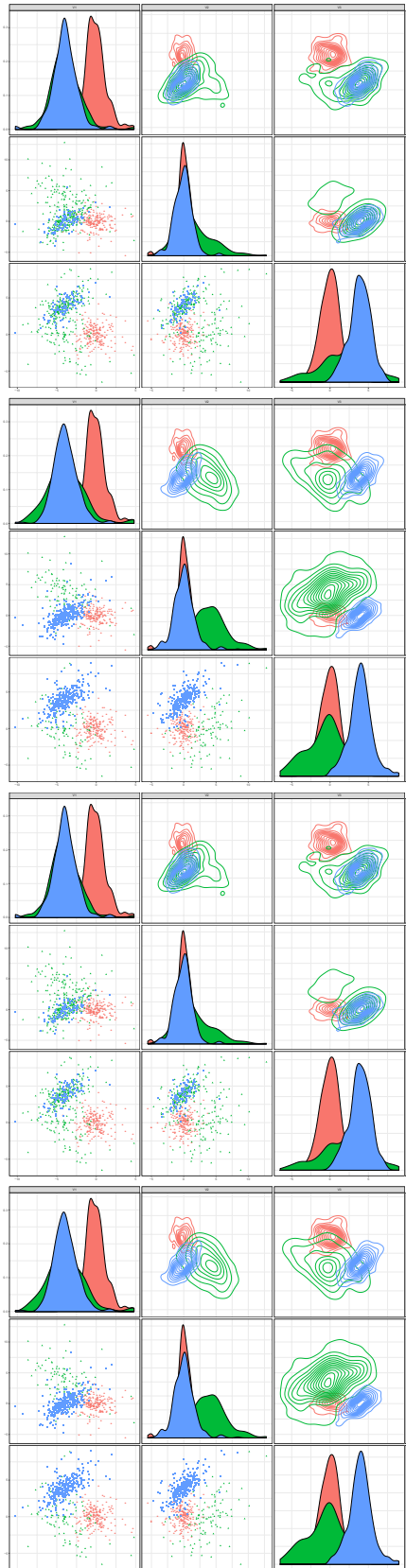


Fig.5 Simulated datasets using multivariate Gaussian distributions Scenarios 1–4, each color and symbol representing a different cluster

Fig.6 Simulated datasets using multivariate Gaussian distributions Scenarios 5–8, each color and symbol representing a different cluster

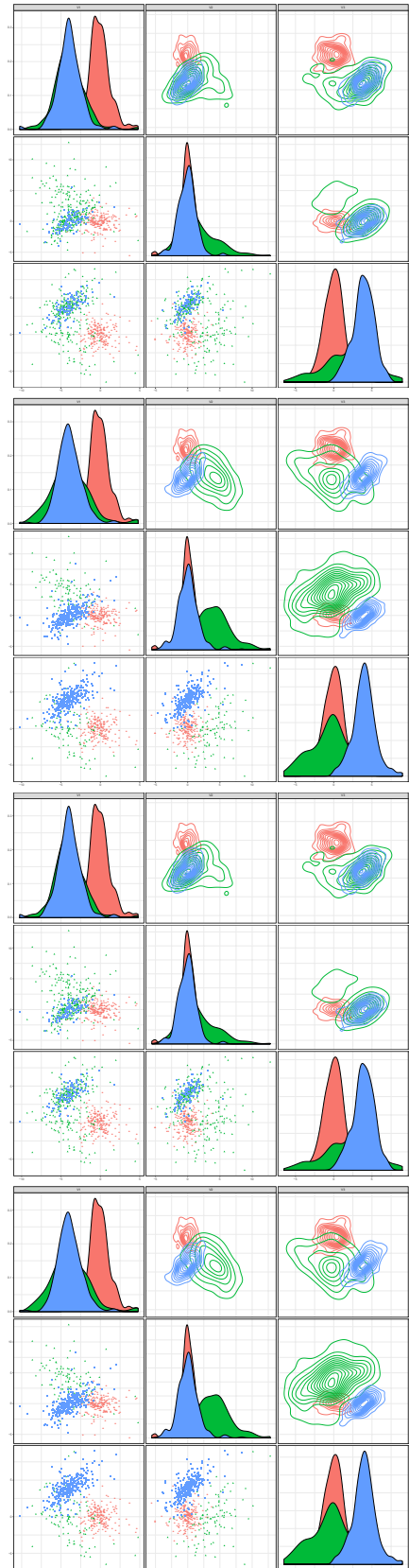
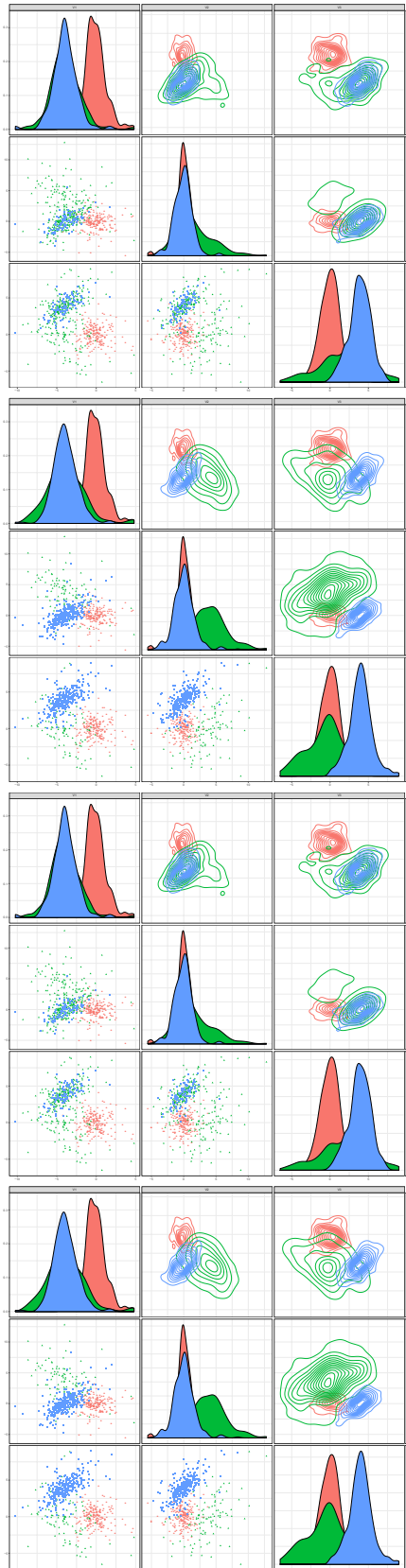


Fig. 7 Simulated datasets using multivariate Gaussian distributions Scenarios 9–12, each color and symbol representing a different cluster

Fig. 8 Simulated datasets using multivariate Gaussian distributions Scenarios 13–16, each color and symbol representing a different cluster

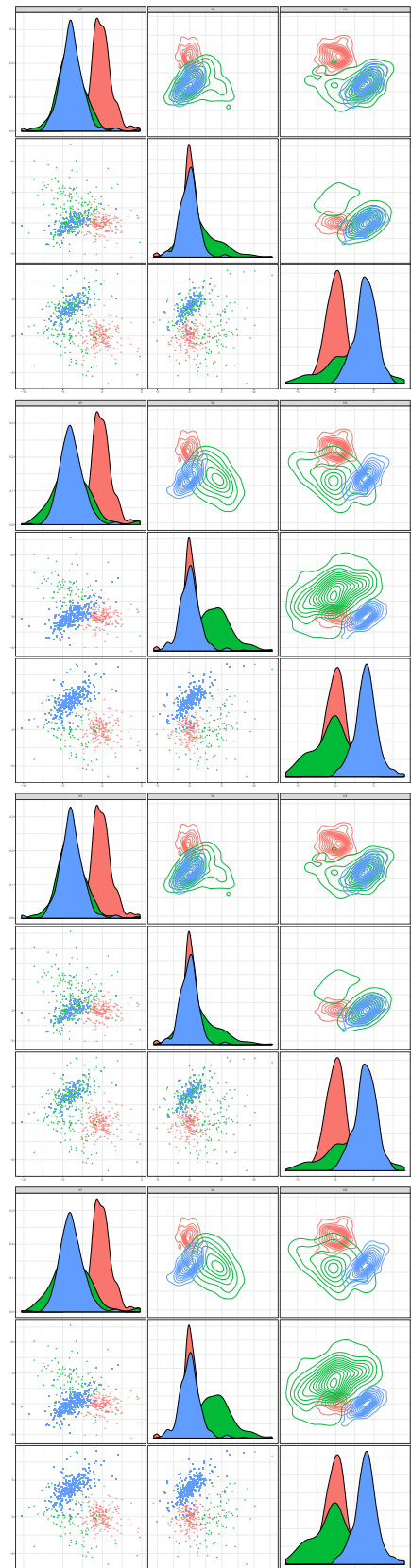
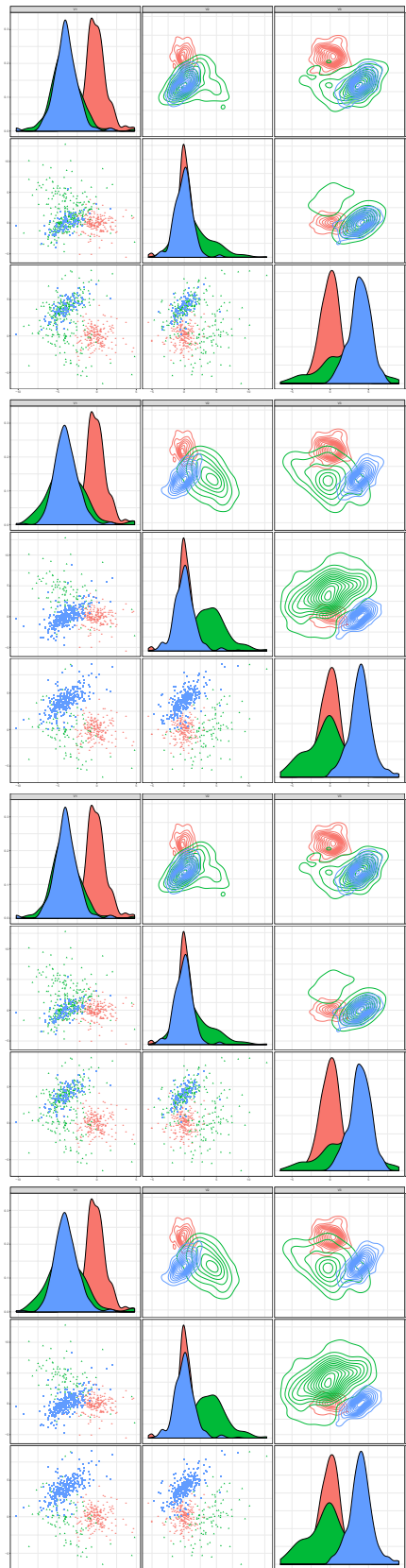


Fig. 9 Simulated datasets using multivariate Student-t distributions Scenarios 1–4, each color and symbol representing a different cluster

Fig. 10 Simulated datasets using multivariate Student-t distributions Scenarios 5–8, each color and symbol representing a different cluster

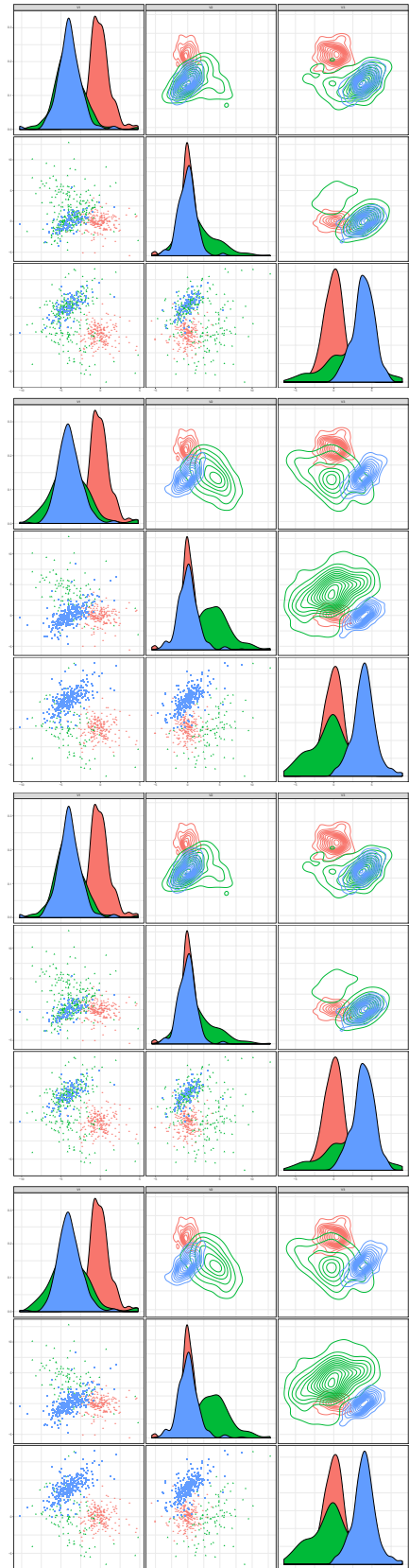
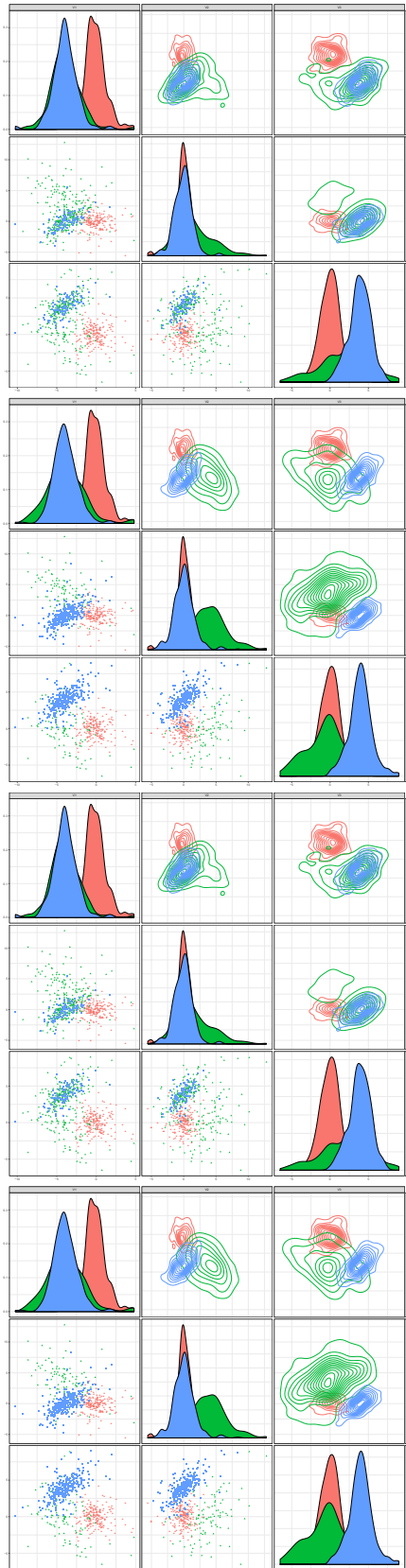


Fig. 11 Simulated datasets using multivariate Student-*t* distributions Scenarios 9–12, each color and symbol representing a different cluster

Fig. 12 Simulated datasets using multivariate Student-*t* distributions Scenarios 13–16, each color and symbol representing a different cluster

468 **References**

469 1. Andrews JL, Wickins JR, Boers NM, McNicholas PDT. An r pack- 519
 470 age for model-based clustering and classification via the multi- 520
 471 variate t distribution. *J Stat Softw.* 2018;83:7. 521
 472 2. Barnett V. Comparative statistical inference. 3rd ed. Hoboken: 522
 473 Wiley; 1999. 523
 474 3. Ben-Israel A, Iyigun C. Probabilistic d-clustering. *J Classif.* 524
 475 2008;25(1):5–26. 525
 476 4. Bezdek JC, Ehrlich R, Full W. Fcm: the fuzzy c-means clustering 526
 477 algorithm. *Comput Geosci.* 1984;10(2–3):191–203. 527
 478 5. Blight B. Estimation from a censored sample for an exponential 528
 479 family. *Biometrika.* 1970;57:389–95. 529
 480 6. Browne RP, ElSherbiny A, McNicholas PD. mixture: mixture 530
 481 models for clustering and classification; R package version 1.4. 531
 482 2015. 532
 483 7. Buck S. A method of estimation of missing values in multivariate 533
 484 data suitable for use with an electronic computer. *J R Stat Soc B.* 534
 485 1960;22:302–6. 535
 486 8. Chiang M, Mirkin B. Intelligent choice of the number of clusters 536
 487 in k-means clustering: an experimental study with different cluster 537
 488 spreads. *J Classif.* 2010;27(1):3–40. 538
 489 9. Dang UJ, Browne RP, McNicholas PD. Mixtures of multivariate 539
 490 power exponential distributions. *Biometrics.* 2015;71(4):1081–9. 540
 491 <https://doi.org/10.1111/biom.12351>. 541
 492 10. Dempster AP, Laird NM, Rubin DB. Maximum likelihood 542
 493 from incomplete data via the EM algorithm. *J R Stat Soc B.* 543
 494 1977;39(1):1–38. 544
 495 11. Genz A, Bretz F, Miwa T, Mi X, Leisch F, Scheipl F, Hothorn 545
 496 T. mvtnorm: multivariate Normal and t Distributions; R package 546
 497 version 1.0-8. 2018. 547
 498 12. Hubert L, Arabie P. Comparing partitions. *J Classif.* 548
 499 1985;2(1):193–218. 549
 500 13. Iyigun C. Probabilistic Distance Clustering. Ph.D. thesis, New 550
 501 Brunswick Rutgers, The State University of New Jersey, 2007. 551
 502 14. Iyigun C, Ben-Israel A. Probabilistic distance clustering adjusted 552
 503 for cluster size. *Prob Eng Inf Sci.* 2008;22(04):603–21. 553
 504 15. Kaufman L, Rousseeuw P. Finding groups in data: an introduction 554
 505 to cluster analysis. New York: Wiley; 1990. 555
 506 16. Kulin HW, Kuenne RE. An efficient algorithm for the 556
 507 numerical solution of the generalized weber problem in spatial 557
 508 economics. *J Reg Sci.* 1962;4(2):21–33. [https://doi.](https://doi.org/10.1111/j.1467-9787.1962.tb00902.x) 558
 509 [org/10.1111/j.1467-9787.1962.tb00902.x](https://doi.org/10.1111/j.1467-9787.1962.tb00902.x). 559
 510 17. Lange KL, Little RJ, Taylor JM. Robust statistical modeling using 560
 511 the t distribution. *J Am Stat Assoc.* 1989;84(408):881–96. 561
 512 18. Lee SX, McLachlan GJ. Finite mixtures of multivariate skew 562
 513 t-distributions: some recent and new results. *Stat Comput.* 563
 514 2014;24(2):181–202. 564
 515 19. Lin TI. Robust mixture modeling using multivariate skew t distri- 565
 516 butions. *Stat Comput.* 2010;20(3):343–56. 566
 517 20. MacQueen J. Some methods for classification and analysis of mul- 567
 518 ti-variate observations. *Proc Fifth Berkeley Symp.* 1967;1:281–97.
 21. McNicholas SM, McNicholas PD, Browne RP. A mixture of 519
 variance-gamma factor analyzers. In: Ahmed SE, editor. Big and 520
 complex data analysis: methodologies and applications. Cham: 521
 Springer International Publishing; 2017. p. 369–85. 522
 22. Murray PM, McNicholas PD, Browne RB. A mixture of common 523
 skew-t factor analyzers. *Statistics.* 2014;3(1):68–82. 524
 23. Newcomb S. A generalized theory of the combination of observa- 525
 tion so as to obtain the best result. *Am J Math.* 1886;8:343–66. 526
 24. Orchard T, Woodbury M. A missing information principle: The- 527
 ory and applications. In: C.U.o.C.P. Berkley (ed.) Proceedings 528
 of the 6th Berkeley Symposium on Mathematical Statistics and 529
 Probability; 1972, vol 1, pp. 697–715 530
 25. Punzo A, McNicholas PD. Parsimonious mixtures of multivariate 531
 contaminated normal distributions. *Biometr J.* 532
 2016;58(6):1506–37. 533
 26. R Core Team: R: A language and environment for statistical com- 534
 puting. R Foundation for Statistical Computing, Vienna, Austria. 535
 2018. 536
 27. R Core Team and contributors worldwide: stats: the R Stats Pack- 537
 age 2014; R package version 3.1.2. 2014. 538
 28. Rachev ST, Klebanov LB, Stoyanov SV, Fabozzi FJ. The methods 539
 of distances in the theory of probability and statistics. Berlin: 540
 Springer; 2013. 541
 29. Rainey C, Tortora C, Palumbo F. A parametric version of proba- 542
 bilistic distance clustering. In: Greselin F, Deldossi L, Vichi M, 543
 Bagnato L, editors. Advances in statistical models for data analy- 544
 sis, studies in classification, data analysis, and knowledge organi- 545
 zation. Cham: Springer; 2019. p. 33–43. 546
 30. Rand WM. Objective criteria for the evaluation of clustering 547
 methods. *J Am Stat Assoc.* 1971;66:846–50. 548
 31. Steinley D. Properties of the Hubert-Arabie adjusted Rand index. 549
Psychol Methods. 2004;9(3):386. 550
 32. Tang Y, Browne RP, McNicholas PD. Flexible clustering of high- 551
 dimensional data via mixtures of joint generalized hyperbolic dis- 552
 tributions. *Statistics.* 2018;7(1):e177. 553
 33. Theodoridis S, Koutroumbas K. Pattern recognition. 2nd ed. New 554
 York: Academic Press; 2003. 555
 34. Tortora C, Franczak BC, Browne RP, McNicholas PD. A mix- 556
 ture of coalesced generalized hyperbolic distributions. *J Classif.* 557
 2019;36(1):26–57. 558
 35. Tortora C, Gettler Summa M, Marino M, Palumbo F. Fac- 559
 tor probabilistic distance clustering (FPDC): a new clustering 560
 method for high dimensional data sets. *Adv Data Anal Classif.* 561
 2016;10(4):441–64. 562
 36. Tortora C, Gettler Summa M, Palumbo F. Factor PD-clustering. 563
 In: Berthold UL, Dirk V (eds) Algorithms from and for Nature 564
 and Life; 2013, p. 115–123. 565
 37. Tortora C, McNicholas PD. FPDclustering: PD-clustering and 566
 factor PD-clustering. R package version 1.4. 2019. 567

Publisher's Note Springer Nature remains neutral with regard to 568
 jurisdictional claims in published maps and institutional affiliations. 569

Author Proof