

A Practical Outlier Detection Approach for Mixed-Attribute Data

Mohamed Bouguessa

*University of Quebec at Montreal
Department of Computer Science
Montreal, Qc, Canada
bouguessa.mohamed@uqam.ca*

Abstract

Outlier detection in mixed-attribute space is a challenging problem for which only few approaches have been proposed. However, such existing methods suffer from the fact that there is a lack of an automatic mechanism to formally discriminate between outliers and inliers. In fact, a common approach to outlier identification is to estimate an outlier score for each object and then provide a ranked list of points, expecting outliers to come first. A major problem of such an approach is where to stop reading the ranked list? How many points should be chosen as outliers? Other methods, instead of outlier ranking, implement various strategies that depend on user-specified thresholds to discriminate outliers from inliers. *Ad hoc* threshold values are often used. With such an unprincipled approach it is impossible to be objective or consistent. To alleviate these problems, we propose a principled approach based on the bivariate beta mixture model to identify outliers in mixed-attribute data. The proposed approach is able to automatically discriminate outliers from inliers and it can be applied to both mixed-type attribute and single-type (numerical or categorical) attribute data without any feature transformation. Our experimental study demonstrates the suitability of the proposed approach in comparison to mainstream methods.

Keywords: Data Mining, Outlier detection, Mixed-attribute data, Mixture model, Bivariate beta.

1 **1. Introduction**

2 Outlier detection is the practice of identifying data points which are consider-
3 ably different from the remaining data (Aggarwal, 2013; Cao, Si, Zhang, and Jia,
4 2010; Kriegel, Kroger, Schubert, and Zimek, 2011; Tan, Steinbach, and Kumar,
5 2006). Outlier detection is also known as exception mining or deviation detec-
6 tion because outlier points are exceptional in some sense or they have attribute
7 values that deviate significantly from the expected or typical attribute values
8 (Tan et al., 2006). Identifying outliers has practical applications in different do-
9 mains such as intrusion and fraud detection, medical diagnosis, and many others
10 (Fustes, Dafonte, Arcay, Manteiga, Smith, Vallenari, and Luri, 2013; Maervoet,
11 Vens, Berghe, Blockeel, and Causmaecker, 2012; Alan and Catal, 2011). For
12 example, in medical diagnosis, outliers may arise when the patient is afflicted
13 with some disease, or suffers side-effects from a drug. Efficient detection of such
14 outliers aids in identifying, preventing, and repairing the effects of malicious or
15 faulty behavior (Penny and Jolliffe, 2011).

16 Approaches to outlier detection can be categorised as supervised, semi-
17 supervised, and unsupervised (Angiulli and Fassetti, 2014). In principle, super-
18 vised, as well as semi-supervised learning methods, use labeled data to create
19 a model which distinguishes outliers from inliers. On the other hand, unsuper-
20 vised approaches do not require any labeled objects and detect outliers as points
21 that are considerably dissimilar or inconsistent with respect to the remaining
22 data using some quantified measures of outlierness (Aggarwal, 2013). To im-
23 plement supervised and semi-supervised outlier detection methods, we should
24 first label the training data (Wu and Wang, 2013). The problem here is that
25 labeled data samples are more difficult, expensive and time consuming to obtain
26 than unlabeled ones. This is why unsupervised approaches are more generally
27 and widely used, since they do not require labeled information. In this paper
28 we focus only on unsupervised outlier detection. For more surveys and details
29 on outlier analysis, we refer the reader to Aggarwal (2013). In the following,
30 we first describe some background information by providing a brief description

31 of the key idea of some outlier detection approaches which are relevant to this
32 work. Next, we discuss a number of elements that motivate this study and
33 describe our contributions.

34 *1.1. Background Information*

35 Several unsupervised approaches have been proposed to identify outliers in
36 numerical data. Such approaches can be broadly classified as statistical-based,
37 distance-based, and density-based (Angiulli and Pizzuti, 2005). Statistical-
38 based approaches attempt to fit the data set under investigation to a certain kind
39 of distribution model (in general, the Gaussian model) (Yamanishi, Takeuchi,
40 Williams, and Milne, 2000). Inliers occur in a high probability region of the
41 model while outliers deviate strongly from the distribution. Distance-based
42 approaches evaluate the outlierness of a point based on the distances to its k -
43 nearest neighbors (k NN) (Angiulli and Pizzuti, 2005, 2002). Points with large
44 k NN distance are defined as outliers. Finally, density-based approaches use the
45 number of points within a specific local region of a data point in order to define
46 local density (Breunig, Kriegel, Ng, and Sander, 2000). The local density val-
47 ues could be then used to measure how isolated a point is with respect to the
48 surrounding objects (Wu and Wang, 2013).

49 The aforementioned approaches were specifically designed for numerical data.
50 However, in several applications, attributes in real data sets are not numerical,
51 but have categorical values. For categorical data sets, distance-based as well
52 as density-based techniques must confront the problem of how to choose the
53 measurement of distance or density (Wu and Wang, 2013). This poses a sig-
54 nificant challenge in terms of generalizing algorithms for numerical data to the
55 categorical domain (Aggarwal, 2013). To address this issue, a number of ap-
56 proaches have been proposed to deal with categorical data (Koufakou, Secretan,
57 and Georgiopoulos, 2011; He, Xu, Huang, and Deng, 2005). Some of these ap-
58 proaches use the concept of frequent itemset mining to estimate an outlying
59 score for each point. Inliers are those points which contain sets of items that
60 co-occur frequently in the data sets, while outliers are likely to be the points

		A_1	A_2	A_3	A_4	A_5	A_6	A_7	A_8
Cluster 1	O_1	0.99	0.08	0.69	0.00	Q	A	A	T
	O_2	0.98	0.08	0.70	0.76	F	A	K	D
	O_3	0.96	0.09	0.71	1.00	G	F	W	W
	O_4	0.96	0.09	0.71	0.21	L	A	C	X
	O_5	1.00	0.10	0.68	0.08	L	A	D	Y
Cluster 2	O_6	0.16	0.08	0.26	0.55	M	N	X	S
	O_7	0.50	0.08	0.00	0.56	M	F	X	I
	O_8	0.70	0.10	0.94	0.58	M	M	X	M
	O_9	0.00	0.09	0.47	0.57	M	J	X	W
	O_{10}	0.32	0.08	0.34	0.56	M	G	H	P
Cluster 3	O_{11}	0.04	1.00	0.11	0.93	T	H	B	H
	O_{12}	0.04	0.32	1.00	0.95	N	T	C	H
	O_{13}	0.05	0.88	0.55	0.93	B	R	B	H
	O_{14}	0.05	0.23	0.17	0.94	O	P	B	H
	O_{15}	0.04	0.00	0.03	0.94	A	W	B	D
Outliers	O_{16}	0.60	0.67	0.09	0.11	C	X	F	Q
	O_{17}	0.73	0.91	0.95	0.96	D	Z	G	F
	O_{18}	0.40	0.74	0.31	0.29	R	U	K	C

Figure 1: Mixed-attribute data set with clustered objects and outliers.

61 that contain rare itemsets (Koufakou et al., 2011).

62 In many cases, categorical and numerical data are found in the same data
63 set, as different attributes. This is referred to as mixed-attribute data (Ag-
64 garwal, 2013). Outliers are those objects containing attribute values that are
65 dissimilar to or inconsistent with the remaining objects in both the numerical
66 and the categorical space (Koufakou and Georgiopoulos, 2010; Otey, Ghoting,
67 and Parthasarathy, 2006). To illustrate, Fig. 1 shows a small data set composed
68 of 18 objects with four numerical attributes (A_1, A_2, A_3 , and A_4) and four cat-
69 egorical attributes (A_5, A_6, A_7 , and A_8). As can be seen from this figure, data
70 objects O_1, O_2, \dots, O_{15} are grouped into three clusters, while the remaining
71 points, that is, O_{16}, O_{17} , and O_{18} , are outliers which could not be located in
72 any cluster. Note that in this figure each cluster is represented by a shade of
73 gray and the unclustered background is white. Clusters thus exist in different
74 subspaces spanned by different attributes. From Fig. 1, we can see that, in
75 contrast to inliers (that is, the clustered objects), outliers contain dissimilar
76 attribute values. In fact, compared to points that belong to clusters, outliers
77 have non-correlated numerical attribute values along the numerical space and
78 infrequent attribute values across the categorical space. On the other hand,

79 from Fig. 1, we can see that objects grouped within clusters contain attribute
80 values that are closely related along a specific subset of dimensions. For ex-
81 ample, objects O_1, O_2, O_3, O_4 , and O_5 , which form cluster 1, contain correlated
82 attribute values along the numerical attributes A_1, A_2, A_3 , and a large number
83 of common categorical attribute values along the categorical attribute A_6 .

84 In practice, when faced with mixed-attribute data, it is common to discretize
85 the numerical attributes and treat all the data as categorical so that categorical
86 outlier detection algorithms can be applied to the entire data set. However, as
87 suggested in Zhang and Jin (2010), discretizing numerical values into several
88 bins could introduce noise or information losses. Improper discretizing thus
89 would hamper the detection performance. To alleviate this problem, only few
90 approaches (Koufakou and Georgiopoulos, 2010; Zhang and Jin, 2010; Otey,
91 Ghoting, and Parthasarathy, 2006), have been proposed to handle outliers in
92 the mixed-attribute space.

93 The approach proposed in Otey et al. (2006) is based on the concept of
94 frequent itemsets to deal with categorical attributes, and the covariance for
95 continuous attributes. Specifically, the authors in Otey et al. (2006) assign to
96 each point an outlier score inversely proportionate to its infrequent itemsets.
97 They also maintain a covariance matrix for each itemset to compute anomaly
98 scores in the continuous attribute space. A point is likely to be an outlier if it
99 contains infrequent categorical sets, or if its continuous values differ from the
100 covariance violation threshold. It is worth noting that the work proposed by
101 Otey et al. (2006) has the merit of being the first outlier detection algorithm
102 for mixed-attribute data.

103 Koufakou and Georgiopoulos (2010) proposed an approach named ODMAD
104 (Outlier Detection for Mixed Attribute Datasets). This algorithm calculates
105 first, for each point in the categorical space, an outlier score which depends on
106 the infrequent subsets contained in that point. Data points with score values less
107 than a user-entered frequency threshold are isolated since they contain highly
108 infrequent categorical values and may thus potentially correspond to outliers.
109 This process results in a reduced data set based on which other outlier scores

110 are calculated for the numerical space using the cosine similarity measure. As
111 described in Koufakou and Georgiopoulos (2010), since minimum cosine simi-
112 larity is 0 and maximum is 1, the data points with similarity close to 0 are more
113 likely to be outliers. Experiments in Koufakou and Georgiopoulos (2010), show
114 that ODMAD is fast and outperforms Otey’s approach.

115 Zhang and Jin (2010) proposed a Pattern based Outlier Detection approach
116 (POD). Patterns in Zhang and Jin (2010) are defined to describe the data objects
117 as well as to capture interactions among different types of attributes. The more
118 an object deviates from these patterns, the higher its outlier score. The authors
119 in Zhang and Jin (2010) use logistic regression to learn patterns. These patterns
120 are then used to estimate outlier scores for objects with mixed attribute. The
121 top n points with the highest score values are declared as outliers. It is important
122 to note that POD is not able to handle categorical values directly. To detect the
123 target patterns, categorical attributes are first mapped into binary attributes.
124 Then, these binary attributes are analyzed together with the original continuous
125 attributes to detect outliers in the mixed-attribute space.

126 *1.2. Motivations and Contributions*

127 The area of outlier detection in mixed-attribute data offers several oppor-
128 tunities for improvement. There are just very few approaches around in the
129 literature so far, yet there are a number of problems still to solve. For instance,
130 the output of POD (Zhang and Jin, 2010) is a ranked list of points that repre-
131 sents the degree of outlierness of each point. The top n points in the list with
132 the highest degree values are considered as outliers. This method encounters a
133 major concern: at which level should this list be cut? Stated otherwise, starting
134 from the first (ranked number one) object, how far should we go in that list? In
135 general, no principled way is suggested on how many points should be selected
136 from a ranked list. In some situations, the top n points are selected solely on the
137 basis of specific knowledge of an application. Unfortunately, prior knowledge
138 about the data under investigation is not always available.

139 Since a ranked list has a particular disadvantage because there is no clear

140 cut-off point of where to stop consulting the results, thresholding has turned
141 out to be important in detecting outliers. For instance, ODMAD (Koufakou
142 and Georgiopoulos, 2010) and the approach proposed by Otey et al. (2006)
143 implement various strategies that depend on user-specified thresholds to detect
144 outliers. In real situations, however, it is rarely possible for users to supply the
145 threshold values accurately. Outlier detection accuracy can thus be seriously
146 reduced if an incorrect threshold value is used. The experiments conducted in
147 Koufakou and Georgiopoulos (2010) on the impact of using various threshold
148 values on the outlier detection accuracy corroborate our claim. Finally, it is
149 worth noting that ODMAD and Otey’s approach depend also on other input
150 parameters such as the minimum support, the maximum length of itemset and
151 the size of a window of categorical and numerical scores. Setting appropriate
152 values of these parameters is not a straightforward task.

153 To alleviate the aforementioned drawbacks of existing approaches for detect-
154 ing outliers in the mixed-attribute space, we propose in this paper a principled
155 approach which is able to automatically identify outliers. In our approach, we
156 first estimate an outlying score, for each object, in the numerical space and
157 another score in the categorical space. Next, we associate to each data point a
158 two dimensional vector containing the estimated scores: one dimension contains
159 the score estimated in the numerical space while the second one contains the
160 outlying score calculated in the categorical space. We assume that, in both
161 spaces, outliers are characterised by high score values. Finally, we propose a
162 statistical framework, based on the bivariate beta mixture, in order to model
163 the estimated outlier score vectors. The goal is to cluster the estimated vectors
164 into several components such that data points associated to the component with
165 the highest score values correspond to outliers.

166 We have used the beta mixture mainly because it permits multiple modes
167 and asymmetry and can thus approximate a wide variety of shapes (Dean and
168 Nugent, 2013; Bouguila and Elguebaly, 2012; Bouguessa, 2012; Ji, Wu, Liu,
169 Wang, and Coombes, 2005), while several other distributions are not able to do
170 so. For example, the standard Gaussian distribution permits symmetric “bell”

171 shape only. However, in many real life applications, the data under investigation
172 is skewed with non-symmetric shapes. In this setting, as observed in Dean and
173 Nugent (2013), and in Boutemedjet, Ziou, and Bouguila (2011), the standard
174 Gaussian distribution may lead to inaccurate modeling (e.g. over estimation of
175 the number of components in the mixture, increase of misclassification errors,
176 etc.). In contrast to several distributions, the beta distribution is more flexible
177 and powerful since it permits multiple symmetric and asymmetric modes, it
178 may be skewed to the right, skewed to left or symmetric (Bouguila, Ziou, and
179 Monga, 2006). This great shape flexibility of the beta distribution provides
180 a better fitting of the outlier score vectors, which leads, in turn, to accurate
181 detection of outliers. Our experimental results corroborate our claim.

182 We summarize the significance of our work as follows:

- 183 1. We view the task of identifying outliers from a mixture modeling perspec-
184 tive, on which we devise a principled approach which is able to formally
185 discriminate between outliers and inliers, while previous works provide
186 only a ranked list of objects expecting outliers to come first.
- 187 2. The proposed method automatically identifies outliers, while existing ap-
188 proaches require human intervention in order to set a detection threshold
189 or to manually define the number of outliers to be identified. Furthermore,
190 our method is general, in the sense that it is not limited to mixed-attribute
191 data and it can be applied to single-type attribute (numerical or categor-
192 ical) data without any feature transformation.
- 193 3. We conducted detailed experiments on several real data sets with mixed-
194 attribute as well as with single-type attribute. The results suggest that
195 the proposed approach achieves competitive results in comparison to main-
196 stream outlier detection algorithms.

197 The rest of this paper is organized as follows. Section 2 describes our ap-
198 proach in detail. An empirical evaluation of the proposed method is given in
199 Section 3. Finally, our conclusion is given in Section 4.

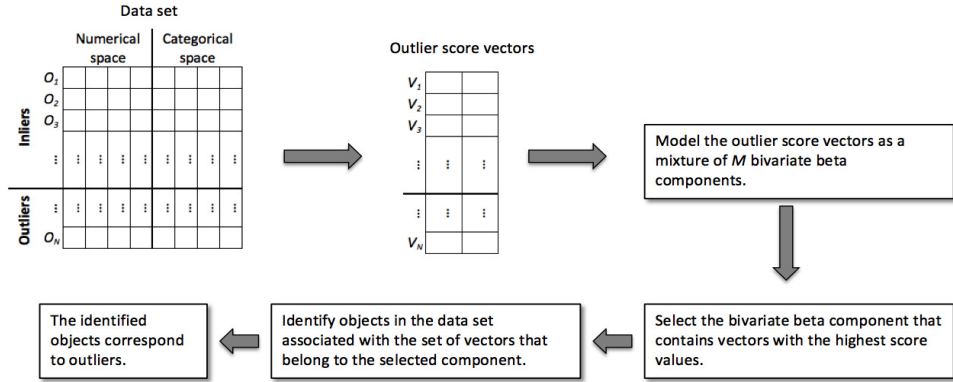


Figure 2: Workflow of the proposed approach.

200 2. Proposed Approach

201 We begin by fixing a proper notation to be used throughout the paper. Let
 202 $\mathcal{D} = \{O_1, \dots, O_N\}$ be a set of N mixed-attribute data points. Each point con-
 203 tains A_n numerical attributes and A_c categorical attributes. The subspace of \mathcal{D}
 204 that contains only numerical attributes is denoted S_{num} , while S_{cat} refers to the
 205 subspace of \mathcal{D} which contains only categorical attributes. In this paper, we rep-
 206 resent a data point O_i as $O_i = [O_i^n, O_i^c]$, such that $O_i^n = (o_{i1}^n, \dots, o_{il}^n, \dots, o_{iA_n}^n)$
 207 and $O_i^c = (o_{i1}^c, \dots, o_{it}^c, \dots, o_{iA_c}^c)$, where o_{il}^n designates the l^{th} numerical attribute
 208 value and o_{it}^c corresponds to the t^{th} categorical attribute value. In what follows,
 209 we will call o_{il}^n a numerical 1D point and o_{it}^c a categorical 1D point.

210 In our approach, we propose first to estimate, for each object O_i , an outlier
 211 score in the numerical space and another score in the categorical space. Then,
 212 we associate to each data point a two-dimensional outlier score vector \vec{V}_i con-
 213 taining the two estimated scores. Finally, based on $\{\vec{V}_i\}_{(i=1, \dots, N)}$, we devise a
 214 probabilistic approach that uses the bivariate beta mixture model to automati-
 215 cally discriminate outliers from inliers in the full-dimensional space. Specifically,
 216 we first model $\{\vec{V}_i\}$ as a mixture of m bivariate beta components. We then se-
 217 lect the component that corresponds to vectors with the highest score values.
 218 Data objects associated with the set of vectors that belong to the selected
 219 component correspond to outliers. Fig. 2 provides a simple visual illustration of the

220 proposed approach. More details are given in the follows.

221 *2.1. Estimating Outlier Score in the Numerical Space*

222 It is widely accepted that outliers are data points that are considerably
 223 dissimilar from the remaining data (Aggarwal, 2013; Huang and Yang, 2011;
 224 Kriegel et al., 2011). In this setting, it is reasonable to assume that, in gen-
 225 eral, most of the attribute values of outlier objects projected along each of the
 226 dimensions in S_{num} tend to be far apart from the remaining attribute values
 227 (Tan et al., 2006). On the other hand, inliers have attribute values that tend to
 228 be closely related along several (or all) dimensions in S_{num} . Our assumption is
 229 based on the fact that inliers tend to form dense regions across several dimen-
 230 sions in the numerical space, while outliers are sparsely distributed. With this
 231 intuition in mind, we define the outlier score $\mathcal{ON}(O_i^n)$ for an object O_i in the
 232 numerical attribute space as

$$\mathcal{ON}(O_i^n) = \sum_{l=1}^{A_n} \log(W_N(o_{il}^n) + 1) \quad (1)$$

233 with

$$W_N(o_{il}^n) = \sum_{j=1}^k \left[d_l(o_{il}^n, kNN_j(o_{il}^n)) \right]^2 \quad (2)$$

234 where, for a specific dimension l in S_{num} , $kNN_j(o_{il}^n)$ denotes the j^{th} nearest (1D
 235 point) neighborhood of o_{il}^n and d_l denotes the distance between two numerical
 236 1D points. In our case, this distance simply corresponds to the absolute value
 237 of the difference between two numerical attribute values of a specific dimension.

238 The outlier score defined in (1) is the sum, over all dimensions in the numer-
 239 ical space S_{num} , of the log of the weight $W_N(o_{il}^n)$. As described by (2), $W_N(o_{il}^n)$
 240 computes the sum of the square of the distance between each 1D point o_{il}^n and
 241 its k nearest neighborhoods in dimension l . Intuitively, a large value of $W_N(o_{il}^n)$
 242 means that o_{il}^n falls into a sparse region in which the k nearest neighborhood
 243 attribute values of o_{il}^n are loosely related, while a small value indicates that o_{il}^n

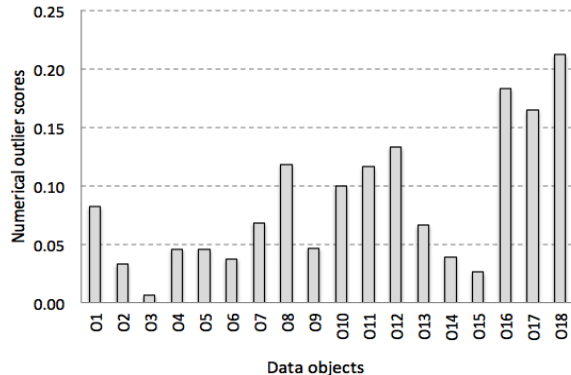


Figure 3: The estimated outlier scores in the numerical space for the data objects depicted by Fig. 1.

244 belongs to a dense region in which the k nearest neighborhood of o_{il}^n are closely
 245 related. Note that we have used the square power in (2) in order to favor the
 246 weight of the 1D points belonging to a sparse region.

247 The weight $W_N(o_{il}^n)$ captures the degree of isolation of an attribute value
 248 with respect to its neighbors. The higher its weight, the more distant are its
 249 neighbors along dimension l of S_{num} . Accordingly, based on (2), we can surmise
 250 that objects that are sparsely distributed over S_{num} will receive high $\mathcal{O}\mathcal{N}(O_i^n)$
 251 values, while related points will receive low score values. This means that out-
 252 liers will be characterized by high score values in contrast to inliers. As an
 253 illustration, Fig. 3 shows the estimated outlier scores in the numerical space
 254 for the data objects depicted by Fig. 1. As can be seen from Fig. 3, outlier
 255 objects O_{16} , O_{17} , and O_{18} have high score values in comparison to inliers, that
 256 is, O_1, \dots, O_{15} .

257 It is important to note that we have used the logarithm function in (1)
 258 primarily to squeeze together the large values that characterize outliers and
 259 stretch out the smallest values, which correspond to inliers. This squeezing and
 260 stretching contributes to enhancing the contrast between largest and smallest
 261 values which helps in distinguishing outliers from the rest of the points. Finally,
 262 note that we have added 1 to $W_N(o_{il}^n)$ in equation (1) to avoid the null value in
 263 the calculation of the logarithm, since it is possible to have $W_N(o_{il}^n) = 0$ in the

264 likely case where an attribute has more than k duplicative values.

265 It is clear that the calculation of the k nearest neighbors is, in general, an
266 expensive task, especially when the number of data points N is very large. How-
267 ever, since we are searching for the k nearest neighbors in the one-dimensional
268 space, we can perform the task in an efficient way by sorting the values in each
269 attribute and limiting the number of distance comparisons to a maximum of
270 $2k$ values. The computation of the k NN distance is sensitive to the value of
271 k , which is a limitation common to all k NN based approaches. However, we
272 believe the problem this limitation creates for our approach does not have a
273 major impact. This is because, since we estimate the k NN distances in the
274 one-dimensional space only, the choice of the value of k is not as critical as in a
275 multi-dimensional case. As suggested in Bouguessa and Wang (2009), to gain a
276 clear idea of the sparseness of the neighborhood of a 1D point, we suggest using
277 $k = \sqrt{N}$ as a default value.

278 2.2. Estimating Outlier Score in the Categorical Space

279 Virtually, as suggested in previous studies (Koufakou et al., 2011; Koufakou
280 and Georgiopoulos, 2010; He et al., 2005), outliers in the categorical space are
281 those points that have infrequent attribute categorical values for all dimensions
282 compared to normal points. This means that every categorical 1D point of
283 outlier objects is infrequent across all dimensions of S_{cat} , while inliers have
284 several categorical 1D points which occur with higher frequency along several (or
285 all) categorical attributes (Koufakou et al., 2011; Koufakou and Georgiopoulos,
286 2010). Based on such a definition, the outlier score $\mathcal{OC}(O_i^c)$ for an object O_i in
287 the categorical attribute space is formulated as

$$\mathcal{OC}(O_i^c) = \sum_{t=1}^{A_c} \log(W_C(o_{it}^c)) \quad (3)$$

288 with

$$W_C(o_{it}^c) = f(o_{it}^c) \quad (4)$$

289 where $f(o_{it}^c)$ denotes the number of times o_{it}^c appears in a specific categorical
 290 dimension t of S_{cat} .

291 $\mathcal{OC}(O_i^c)$ is defined as the sum, across all dimensions in the categorical space
 292 S_{cat} , of the log of the weight $W_C(o_{it}^c)$, which, in turn, corresponds to the occur-
 293 rence frequency of o_{it}^c in the categorical attribute t . Here, it is clear that rare
 294 categorical attribute values projected along dimension t will receive low weight
 295 values, while larger $W_C(o_{it}^c)$ values indicate that o_{it}^c is shared by several objects
 296 within dimension t . Accordingly, based on (3), points that share common cate-
 297 gorical values across S_{cat} will get large $\mathcal{OC}(O_i^c)$ values, while data objects that
 298 have infrequent categorical values across S_{cat} will receive low $\mathcal{OC}(O_i^c)$ values. As
 299 a result, since outliers are those points whose attribute categorical values occur
 300 very rarely along each dimension in S_{cat} (Koufakou et al., 2011), it is easy to
 301 see that small values of $\mathcal{OC}(O_i^c)$ designate outliers and high scores correspond
 302 to inliers. Finally, note that, as with the numerical outlier score described by
 303 (1), we have used a logarithm function in (3) to enhance the contrast between
 304 larger and smaller weight values.

305 In this paper, as mentioned in Section 1, we assume that outliers are char-
 306 acterized by large score values in contrast to inliers. However, as just discussed,
 307 large $\mathcal{OC}(O_i^c)$ scores refer to inliers. To regularize such scores, we need to invert
 308 them. For this purpose, we simply take the difference between the observed
 309 score and the maximum possible estimated score \mathcal{OC}_{max} . The inverted score is
 310 estimated as

$$\mathcal{OC}_{inv}(O_i^c) = \mathcal{OC}_{max} - \mathcal{OC}(O_i^c) \quad (5)$$

311 It easy to show that this linear inversion doesn't affect the ranking-stability of
 312 the inverted scores:

$$\begin{aligned} \mathcal{OC}(O_1^c) \leq \mathcal{OC}(O_2^c) &\iff -\mathcal{OC}(O_1^c) \geq -\mathcal{OC}(O_2^c) \\ &\iff \mathcal{OC}_{max} - \mathcal{OC}(O_1^c) \geq \mathcal{OC}_{max} - \mathcal{OC}(O_2^c) \\ &\iff \mathcal{OC}_{inv}(O_1^c) \geq \mathcal{OC}_{inv}(O_2^c). \end{aligned}$$

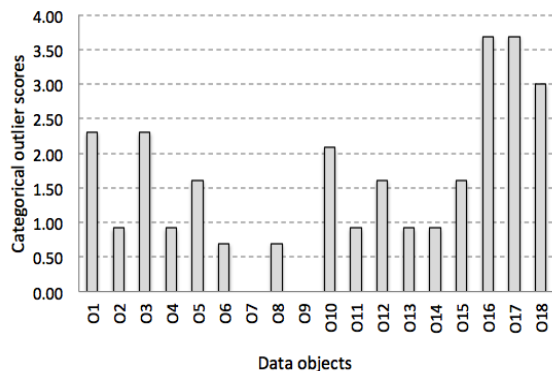


Figure 4: The estimated outlier scores in the categorical space for the data objects depicted by Fig. 1.

313 Accordingly, based on such a linear inversion, outliers will receive large score
 314 values while inliers will receive the lowest score values. In the remainder of this
 315 paper, unless otherwise specified, we use only the inverted categorical outlier
 316 score values. As an illustration, Fig. 4 shows the estimated outlier scores in the
 317 categorical space for the data objects depicted by Fig. 1. As can be seen from
 318 Fig. 4, outlier objects O_{16} , O_{17} , and O_{18} have high score values in comparison
 319 to inliers, that is, O_1, \dots, O_{15} .

320 Finally, as the reader can notice, in our approach we treat numerical and
 321 categorical attributes independently in order to estimate outlier scores in the
 322 numerical and the categorical space. In other words, this means we assume the
 323 independence of both numerical and categorical attributes. Such an assump-
 324 tion is mainly based on the general definition of outliers, which relies on the fact
 325 that outlier objects contain attribute values that are dissimilar to or inconsis-
 326 tent with the remaining points. Stated otherwise, outliers may contain many atyp-
 327 ical attribute values across most (or all) attributes of the data in comparison
 328 to inliers. Accordingly, investigating individual attributes in order to localize
 329 attribute values that deviate significantly from the expected or typical attribute
 330 values is appropriate to effectively detect outliers in the whole space.

331 *2.3. Modeling Outlier Score Vectors*

332 Once the outlier scores are estimated in both the numerical and the categor-
 333 ical spaces, we now focus on how to automatically identify outliers in the mixed-
 334 attribute space. To this end, we associate to each object O_i a two-dimensional
 335 vector \vec{V}_i such that the first element of this vector corresponds to the outlier
 336 score of O_i in the numerical space, while the second element represent the out-
 337 lier score of O_i in the categorical space. Then, based on the estimated vectors,
 338 we propose a probabilistic approach that uses the bivariate beta mixture model
 339 to automatically discriminate outliers from inliers in the full-dimensional space.
 340 The probabilistic model framework is described in the follows.

341 *2.3.1. The Bivariate Beta Mixture Model*

342 Since the beta distribution is defined on the interval $[0,1]$, we should first,
 343 without loss of generality, normalize the estimated outlier score values between
 344 0 and 1. Let $\vec{V}_i = (V_{i1}, V_{i2})^T$ where V_{i1} and V_{i2} represent, respectively, the
 345 normalized outlier scores $\mathcal{ON}(O_i^n)$ and $\mathcal{OC}_{inv}(O_i^c)$. Under a mixture of bivariate
 346 beta distribution,

$$\vec{V}_i \sim \sum_{m=1}^M \lambda_m \mathcal{B}_m(\vec{V}_i | \vec{x}_m, \vec{y}_m) \quad (6)$$

347 where $\mathcal{B}_m(\vec{V}_i | \vec{x}_m, \vec{y}_m)$ is the m^{th} bivariate beta distribution; M denotes the
 348 number of components in the mixture; $\vec{x} = \{\vec{x}_1, \dots, \vec{x}_M\}$ and $\vec{y} = \{\vec{y}_1, \dots, \vec{y}_M\}$.
 349 \vec{x}_m and \vec{y}_m are the parameters of the m^{th} component with $\vec{x}_m = (x_{m1}, x_{m2})^T$
 350 and $\vec{y}_m = (y_{m1}, y_{m2})^T$. $\lambda = \{\lambda_1, \dots, \lambda_M\}$ represents the mixing coefficients
 351 such that $\sum_{m=1}^M \lambda_m = 1$ and $\lambda_m > 0$.

352 The bivariate beta distribution can be obtained by cascading two beta vari-
 353 ables together, that is, each element in the two-dimensional vector \vec{V}_i is a scalar
 354 beta variable. In other words, the bivariate beta is the product of two uni-
 355 variate beta densities. Accordingly, the probability density function of the m^{th}

356 bivariate beta component is expressed as

$$\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m) = \prod_{d=1}^2 \mathcal{B}(V_{id}|x_{md}, y_{md}) \quad (7)$$

357 $\mathcal{B}(V_{id}|x_{md}, y_{md})$ is the probability density function of the univariate beta distri-
 358 bution which is given by

$$\mathcal{B}(V_{id}|x_{md}, y_{md}) = \frac{\Gamma(x_{md} + y_{md})}{\Gamma(x_{md})\Gamma(y_{md})} V_{id}^{x_{md}-1} (1 - V_{id})^{y_{md}-1} \quad (8)$$

359 where $\Gamma(\cdot)$ is the gamma function given by $\Gamma(\alpha) = \int_0^\infty \beta^{\alpha-1} \exp(-\beta) d\beta; \beta > 0$.

360 2.3.2. Maximum Likelihood Estimate

361 A common approach for estimating the unknown parameters x_{md} and y_{md} ,
 362 ($m = 1, \dots, M; d = 1, 2$) is the maximum likelihood estimation technique. The
 363 likelihood function corresponding to the m^{th} bivariate beta component \mathcal{B}_m is
 364 defined as

$$\begin{aligned} \mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)) &= \prod_{\vec{V}_i \in \mathcal{B}_m} \mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m) \\ &= \prod_{\vec{V}_i \in \mathcal{B}_m} \prod_{d=1}^2 \mathcal{B}(V_{id}|x_{md}, y_{md}) \end{aligned} \quad (9)$$

365 The logarithm of the likelihood function is given by

$$\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m))] = \sum_{i=1}^{N_m} \sum_{d=1}^2 \log [\mathcal{B}(V_{id}|x_{md}, y_{md})] \quad (10)$$

366 where N_m is the size of the m^{th} component.

We note that the parameters pair $\{x_{md}, y_{md}\}$ is independent from all other pairs. The problem of estimating the parameters of the model can thus be reduced to the estimation of the parameters pair $\{x_{md}, y_{md}\}$ independently over each dimension of the outlier score vectors belonging to component m . In this

setting, the value $\{\widehat{x}_{md}, \widehat{y}_{md}\}$ that maximizes the likelihood can be obtained by taking the derivative of the expectation of the log-likelihood function with respect to x_{md} and y_{md} and setting the gradient equal to zero as

$$\begin{bmatrix} \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)])]}{\partial x_{md}} \\ \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)])]}{\partial y_{md}} \end{bmatrix} = 0 \quad (11)$$

367 where

$$\begin{aligned} \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)])]}{\partial x_{md}} &= \sum_{i=1}^{N_m} \left[\frac{\partial}{\partial x_{md}} \log \left(\frac{\Gamma(x_{md} + y_{md})}{\Gamma(x_{md})\Gamma(y_{md})} V_{id}^{x_{md}-1} (1 - V_{id})^{y_{md}-1} \right) \right] \\ &= \sum_{i=1}^{N_m} \left[\frac{\Gamma'(x_{md} + y_{md})}{\Gamma(x_{md} + y_{md})} - \frac{\Gamma'(x_{md})}{\Gamma(x_{md})} + \log(V_{id}) \right] \\ &= N_m \frac{\Gamma'(x_{md} + y_{md})}{\Gamma(x_{md} + y_{md})} - N_m \frac{\Gamma'(x_{md})}{\Gamma(x_{md})} + \sum_{i=1}^{N_m} \log(V_{id}). \end{aligned} \quad (12)$$

368 and

$$\begin{aligned} \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)])]}{\partial y_{md}} &= \sum_{i=1}^{N_m} \left[\frac{\partial}{\partial y_{md}} \log \left(\frac{\Gamma(x_{md} + y_{md})}{\Gamma(x_{md})\Gamma(y_{md})} V_{id}^{x_{md}-1} (1 - V_{id})^{y_{md}-1} \right) \right] \\ &= \sum_{i=1}^{N_m} \left[\frac{\Gamma'(x_{md} + y_{md})}{\Gamma(x_{md} + y_{md})} - \frac{\Gamma'(y_{md})}{\Gamma(y_{md})} + \log(1 - V_{id}) \right] \\ &= N_m \frac{\Gamma'(x_{md} + y_{md})}{\Gamma(x_{md} + y_{md})} - N_m \frac{\Gamma'(y_{md})}{\Gamma(y_{md})} + \sum_{i=1}^{N_m} \log(1 - V_{id}). \end{aligned} \quad (13)$$

369 Equations (11), (12) and (13) yield the following expression

$$\begin{bmatrix} N_m [\psi(x_{md} + y_{md}) - \psi(x_{md})] + \sum_{i=1}^{N_m} \log(V_{id}) \\ N_m [\psi(x_{md} + y_{md}) - \psi(y_{md})] + \sum_{i=1}^{N_m} \log(1 - V_{id}) \end{bmatrix} = 0 \quad (14)$$

370

371 where $\psi(\cdot)$ is the digamma function given by $\psi(\alpha) = \frac{\Gamma'(\alpha)}{\Gamma(\alpha)}$.

372 Since the digamma function is defined through an integration, a closed-
 373 form solution to (14) does not exist. So the parameters pair $\{x_{md}, y_{md}\}$ can
 374 be estimated using the Newton-Raphson method (Ypma, 1995). Specifically,
 375 $\{x_{md}, y_{md}\}$ are estimated iteratively:

$$\begin{bmatrix} x_{md}^{(I+1)} \\ y_{md}^{(I+1)} \end{bmatrix} = \begin{bmatrix} x_{md}^{(I)} \\ y_{md}^{(I)} \end{bmatrix} - [\vec{h}_m]^T [\mathcal{H}_m]^{-1} \quad (15)$$

376 where I is the iteration index, h_m and \mathcal{H}_m are respectively the vector of the
 377 first derivatives and the matrix of the second derivatives of the log likelihood
 378 function of the m^{th} component.

379 The vector \vec{h}_m is defined as

$$\vec{h}_m = \begin{pmatrix} h_m^1 \\ h_m^2 \end{pmatrix} = \begin{pmatrix} \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i | x_m^-, y_m^-)])]}{\partial x_{md}} \\ \frac{\partial E(\log [\mathcal{L}(\mathcal{B}_m(\vec{V}_i | x_m^-, y_m^-)])]}{\partial y_{md}} \end{pmatrix} \quad (16)$$

380 and the matrix \mathcal{H}_m is expressed as

$$\mathcal{H}_m = \begin{pmatrix} \frac{\partial h_m^1}{\partial x_{md}} & \frac{\partial h_m^1}{\partial y_{md}} \\ \frac{\partial h_m^2}{\partial x_{md}} & \frac{\partial h_m^2}{\partial y_{md}} \end{pmatrix}, \quad (17)$$

382 where

$$\begin{aligned} \frac{\partial h_m^1}{\partial x_m} &= N_m [\psi'(x_{md} + y_{md}) - \psi'(x_{md})], \\ \frac{\partial h_m^1}{\partial y_m} &= \frac{\partial h_m^2}{\partial x_{md}} = N_m [\psi'(x_{md} + y_{md})], \\ \frac{\partial h_m^2}{\partial y_m} &= N_m [\psi'(x_{md} + y_{md}) - \psi'(\beta_{md})]. \end{aligned} \quad (17)$$

383 $\psi'(\cdot)$ is the trigamma function given by $\psi'(\alpha) = \frac{\Gamma''(\alpha)}{\Gamma(\alpha)} - [\frac{\Gamma'(\alpha)}{\Gamma(\alpha)}]^2$.

384 The Newton-Raphson algorithm for the update of equation (15) converges,
 385 as our estimate of x_{md} and y_{md} change by less than a small positive value ϵ
 386 with each successive iteration, to \hat{x}_{md} and \hat{y}_{md} . Note that, we have used in

387 our implementation the method of moments estimators of the beta distribution
388 (Bain and Engelhardt, 2000) to define starting values for $\{x_{md}^{(0)}, y_{md}^{(0)}\}$ in equation
389 (15). In this technique, the expected mean of the distribution is equated to the
390 sample mean and the expected variance to the sample variance. Specifically,
391 the method of moments estimators are

$$\begin{aligned}\hat{x}_{md}^{(0)} &= \bar{\mu}_{md} \left[\frac{\bar{\mu}_{md}(1 - \bar{\mu}_{md})}{\sigma_{md}^2} - 1 \right], \\ \hat{y}_{md}^{(0)} &= (1 - \bar{\mu}_{md}) \left[\frac{\bar{\mu}_{md}(1 - \bar{\mu}_{md})}{\sigma_{md}^2} - 1 \right].\end{aligned}\quad (18)$$

392 where $\bar{\mu}_{md}$ and σ_{md}^2 denote respectively the sample mean and variance of the
393 normalized outlier score vectors belonging to the m^{th} component which are
394 projected along dimension d .

395 2.3.3. EM Algorithm for the Bivariate Beta Mixture

396 Let $\mathcal{P} = \{\lambda_1, \dots, \lambda_M, \vec{x}_1, \dots, \vec{x}_M, \vec{y}_1, \dots, \vec{y}_M\}$ denote the set of parameters
397 of the mixture and $\mathcal{V} = \{\vec{V}_1, \dots, \vec{V}_N\}$ the set of the normalized outlier score
398 vectors. The usual choice for obtaining the maximum likelihood of the distribu-
399 tion parameters is the EM algorithm (Dempster, Laird, and Rubin, 1977). This
400 algorithm is based on the interpretation of \mathcal{V} as *incomplete* data. As mentioned
401 in Figueiredo and Jain (2002), for finite mixture, the *missing* part is a set of N
402 label vectors $\eta = \{\vec{\eta}_1, \dots, \vec{\eta}_N\}$ associated with the N outlier score vectors, in-
403 dicating to which component \vec{V}_i belongs. Specifically, each $\vec{\eta}_i = (\eta_{i1}, \dots, \eta_{im})^T$
404 is a binary vector, where $\eta_{im} = 1$ if \vec{V}_i belongs to component m and $\eta_{im} = 0$
405 otherwise.

406 The complete data is thus defined by the sets η and \mathcal{V} . The likelihood of the
407 complete data is then:

$$\mathcal{L}(\mathcal{V}, \eta | \mathcal{P}) = \prod_{i=1}^N \prod_{m=1}^M [\lambda_m \mathcal{B}_m(\vec{V}_i | \vec{x}_m, \vec{y}_m)]^{\eta_{im}} \quad (19)$$

408 and the complete log likelihood is:

$$\begin{aligned}
\log(\mathcal{L}(\mathcal{V}, \eta|\mathcal{P})) &= \sum_{i=1}^N \sum_{m=1}^M \eta_{im} \log [\lambda_m \mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)] \\
&= \sum_{i=1}^N \sum_{m=1}^M \eta_{im} \log [\lambda_m \prod_{d=1}^2 \mathcal{B}(V_{id}|x_{md}, y_{md})] \\
&= \sum_{i=1}^N \sum_{m=1}^M \eta_{im} [\log(\lambda_m) + \sum_{d=1}^2 \log(\mathcal{B}(V_{id}|x_{md}, y_{md}))] \quad (20)
\end{aligned}$$

409 The EM algorithm can now be used to estimate \mathcal{P} . Specifically, the algo-
410 rithm iterates between an Expectation step and an Maximization step in order
411 to produce a sequence estimate $\{\hat{\mathcal{P}}\}^{(I)}$, ($I = 0, 1, 2, \dots$), where I denotes the
412 current iteration step, until the change in the value of the complete log-likelihood
413 in (20) is negligible. Details of each step are given below.

414 In the Expectation step: each latent variable η_{im} is replaced by its expecta-
415 tion as follows

$$\hat{\eta}_{im}^{(I)} = E[\eta_{im}|\vec{V}_i, \mathcal{P}] = \frac{\hat{\lambda}_m^{(I)} \mathcal{B}_m(\vec{V}_i|\vec{x}_m, \vec{y}_m)}{\sum_{j=1}^M \hat{\lambda}_j^{(I)} \mathcal{B}_j(\vec{V}_i|\vec{x}_j, \vec{y}_j)} \quad (21)$$

416 In the Maximization step: the mixing coefficients $\{\lambda_m\}$ and the parameters
417 $\{\vec{x}_1, \dots, \vec{x}_M, \vec{y}_1, \dots, \vec{y}_M\}$ are calculated using the values of $\hat{\eta}_{im}$ estimated in the
418 Expectation step. Specifically, the mixing coefficients are calculated as

$$\hat{\lambda}_m^{(I+1)} = \frac{\sum_{i=1}^N \hat{\eta}_{im}^{(I)}}{N}, \quad m = 1, \dots, M \quad (22)$$

419 The parameters $\{\vec{x}_m = (x_{m1}, x_{m2})^T\}_{(m=1, \dots, M)}$ and $\{\vec{y}_m = (y_{m1}, y_{m2})^T\}_{(m=1, \dots, M)}$
420 are estimated using the Newton-Raphson algorithm, based on (15), as described
421 in the previous subsection.

422 Finally, note that, the EM algorithm requires the initial parameters of each
423 component. Since EM is highly dependent on initialization, it will be help-
424 ful to perform initialization by means of clustering algorithms (Figueiredo and
425 Jain, 2002). For this purpose we implement the k-means algorithm in order to

Algorithm 1: EM algorithm for the bivariate beta mixture

Input : $\{\vec{V}_i\}_{(i=1,\dots,N)}$; M
Output: $\hat{\mathcal{P}} = \{\hat{\lambda}_1, \dots, \hat{\lambda}_M, \vec{\hat{x}}_1, \dots, \vec{\hat{x}}_M, \vec{\hat{y}}_1, \dots, \vec{\hat{y}}_M\}$
1 begin
 // Initialization
2 Apply the k-means algorithm to cluster the set $\{\vec{V}_i\}$ into M components;
3 Estimate the initial set of parameters of each component using (18);
4 **repeat**
 // Expectation
5 Estimate $\{\hat{\eta}_{im}\}_{(i=1,\dots,N; m=1,\dots,M)}$ using (21);
 // Maximization
6 Estimate $\{\hat{\lambda}_m\}_{(m=1,\dots,M)}$ using (22);
7 Estimate $\{\vec{\hat{x}}_{md}, \vec{\hat{y}}_{md}\}_{(m=1,\dots,M; d=1,2)}$ using (15);
8 **until** the change in (20) is negligible;
9 Return $\hat{\mathcal{P}}$;
10 end

426 partition the set $\{\vec{V}_i\}_{(i=1,\dots,N)}$, into M components. Then, based on such partition,
 427 we estimate the initial parameters of each component using the method
 428 of moment estimator of the beta distribution (Bain and Engelhardt, 2000) and
 429 set them as initial parameters to the EM algorithm. The detailed algorithm is
 430 summarized in Algorithm 1.

431 2.3.4. Estimating the Optimal Number of Components in the Mixture

432 The use of mixture of the bivariate beta distribution allows us to give a
 433 flexible model to describe the outlier score vectors. To form such a model,
 434 we need to estimate M , the number of components, and the parameters for
 435 each component. Several model selection approaches have been proposed to
 436 estimate M (Bouguessa, Wang, and Sun, 2006). In this paper, we implemented
 437 a deterministic approach that uses the EM algorithm described in Algorithm
 438 1 in order to obtain a set of candidate models for the range value of M (from
 439 1 to M_{max} , the maximum number of components in the mixture) which is
 440 assumed to contain the optimal M (Figueiredo and Jain, 2002). The number of
 441 components is then selected according to

$$\hat{M} = \underset{M}{\operatorname{argmin}} \left\{ \mathcal{C}(\hat{\mathcal{P}}, M) \right\}_{M=1, \dots, M_{max}} \quad (23)$$

Algorithm 2: Estimating the number of components in the mixture

Input : $\{\vec{V}_i\}_{(i=1,\dots,N)}$, M_max
Output: The optimal number of components \hat{M}

```

1 begin
2   for  $M = 1$  to  $M\_max$  do
3     if  $M==1$  then
4       Estimate  $\{\hat{x}_d, \hat{y}_d\}_{d=1,2}$  using (15);
5       Estimate  $ICL - BIC(\hat{\mathcal{P}}, M)$  using (24);
6     else
7       Estimate the parameters of the mixture using Algorithm 1;
8       Estimate  $ICL - BIC(\hat{\mathcal{P}}, M)$  using (24);
9     end
10  end
11  Select  $\hat{M}$ , such that  $\hat{M} = \underset{M}{\operatorname{argmin}} ICL - BIC(\hat{\mathcal{P}}, M)$ ;
12 end

```

where $\mathcal{C}(\hat{\mathcal{P}}, M)$ is some model selection criterion. Ji et al. (2005) found that the Integrated Classification Likelihood-Bayesian Information Criterion (ICL-BIC) performs well in selecting the number of components in the beta mixture. ICL-BIC has been also used in Dean and Nugent (2013) to select the number of beta mixture components. Accordingly, we use in our method ICL-BIC to identify the optimal number of components. The ICL-BIC criterion is given by

$$ICL - BIC(\hat{\mathcal{P}}, M) = -2 \log(\mathcal{L}(\mathcal{V}, \hat{\eta} | \hat{\mathcal{P}})) + \mathcal{Q}_M \log(N) - 2 \sum_{i=1}^N \sum_{m=1}^M \hat{\eta}_{im} \log(\hat{\eta}_{im}) \quad (24)$$

442 where \mathcal{Q}_M denotes the number of parameters of the model with M components
 443 and $\log(\mathcal{L}(\mathcal{V}, \hat{\eta} | \hat{\mathcal{P}}))$ corresponds to logarithm of the likelihood at the maximum
 444 likelihood solution for the investigated mixture model. The number of compo-
 445 nents that minimize $ICL - BIC(\hat{\mathcal{P}}, M)$ is considered to be the optimal value for
 446 M . The procedure for estimating the number of components in the mixture is
 447 summarized in Algorithm 2.

448 2.3.5. Automatic Identification of Outlier

449 Once the optimal number of components has been identified, we focus now
 450 on detecting the bivariate beta component that corresponds to outliers. To this

451 end, we used the results of the EM algorithm in order to derive a classifica-
 452 tion decision about which outlier score vector \vec{V}_i belongs to which component
 453 in the mixture. In fact, the EM algorithm yields the final estimated posterior
 454 probability $\hat{\eta}_{im}$, the value of which represents the posterior probability that \vec{V}_i
 455 belongs to component m . We assign \vec{V}_i to the component that corresponds to
 456 the maximum value of $\hat{\eta}_{im}$. We thus divide the set of outlier score vectors into
 457 several components. As discussed earlier, in our approach we assume that out-
 458 lier points are characterized by high score values. Therefore, we are interested
 459 by the bivariate beta component which contains vectors with the highest score
 460 values. To identify such a component, we first compute, for each component
 461 in the mixture, the average value of the numerical outlier scores and also the
 462 average value of the categorical outlier scores (that is, we compute the average
 463 of V_{i1} and V_{i2} per component). Then, we select the component with the largest
 464 average values as our target component. This simple strategy for determining
 465 which component to pick works well in practice since it fits our assumption,
 466 which is based on the fact that outlier points are characterized by large score
 467 values in both numerical and categorical space. Finally, once our target compo-
 468 nent is identified, we focus on the problem of detecting outlier objects. To this
 469 end, we identify the set of data objects that are associated with the outlier score
 470 vectors \vec{V}_i that belong to the selected component. The identified objects are out-
 471 liers. The steps described in Algorithm 3 can be implemented to automatically
 472 identify outliers.

473 Finally, it is worth noting that the proposed methodology could be also
 474 used to identify outlier objects in single-type (categorical or numerical) at-
 475 tribute data. In this particular case, we propose to associate to each object
 476 only one score ($\mathcal{ON}(O_i^n)$ or $\mathcal{OC}_{inv}(O_i^c)$, depending on the attribute type of
 477 the data under investigation). Then, to automatically discriminate between
 478 outliers and inliers, we can model the estimated scores as a finite mixture dis-
 479 tribution using the univariate beta which is given by (8). Here, the problem
 480 is thus reduced from modeling a set of two-dimensional outlier score vectors
 481 $\{\vec{V}_i\}_{(i=1,\dots,N)}$ (in the case of mixed-attribute data) to modeling a list of scalar

Algorithm 3: Automatic identification of outliers

Input : A data set \mathcal{D}

Output: A set of outliers \mathcal{OUT}

```
1 begin
2   Estimate  $\{\mathcal{ON}(O_i^n)\}_{(i=1,\dots,N)}$  using (1);
3   Estimate  $\{\mathcal{OC}_{inv}(O_i^c)\}_{(i=1,\dots,N)}$  using (3) and (5);
4   Associate, to each object  $O_i$  in  $\mathcal{D}$ , a vector  $\vec{V}_i = (V_{i1}, V_{i2})^T$  where  $V_{i1}$  and
    $V_{i2}$  represent, respectively, the normalized values of  $\mathcal{ON}(O_i^n)$  and
    $\mathcal{OC}_{inv}(O_i^c)$  in  $[0,1]$  ;
5   Apply Algorithm 2 to cluster  $\{\vec{V}_i\}_{(i=1,\dots,N)}$  into  $M$  bivariate beta
   components;
6   Use the results of the EM algorithm to decide about the membership of the
   outlier score vectore  $\vec{V}_i$  in each component;
7   Select the bivariate beta component that contains vectors with the highest
   score values;
8   Identify objects in  $\mathcal{D}$  associated with the set of  $\vec{V}_i$  that belong to the
   selected component and store them in  $\mathcal{OUT}$ ;
9   Return  $\mathcal{OUT}$ ;
10 end
```

482 outlier score values ($\{\mathcal{ON}(O_i^n)\}_{(i=1,\dots,N)}$ or $\{\mathcal{OC}_{inv}(O_i^c)\}_{(i=1,\dots,N)}$). In this set-
483 ting, the parameters of the univariate beta mixture model to be estimated are
484 $\{\lambda_m, x_m, y_m\}_{(m=1,\dots,M)}$. These parameters and the optimal number of compo-
485 nents in the mixture are estimated using the EM algorithm with the Newton-
486 Raphson method and ICL-BIC as described in the above subsections. By doing
487 so, we divide the outlier scores into several populations so that the larger scores
488 can be identified and the associated objects are then declared as outliers.

489 3. Experimental Evaluation

490 In this section, we devise an empirical study to evaluate the suitability of the
491 proposed approach. In the following, we first describe the technique that we have
492 adopted to produce data for use in outlier detection and the performance metrics
493 used in the evaluation. Next, we illustrate the effectiveness of our approach to
494 identify outliers in mixed-attribute data. Finally, we devise further experiments
495 to evaluate the performance of the proposed methodology in detecting outliers
496 in single-type attribute data.

497 *3.1. Data Preparation and Metrics*

498 We draw the attention of the reader to the fact that, at the time of writing
499 this paper, there is a shortage of standard benchmark data that can be used
500 to evaluate outlier detection algorithms. Most of the publicly available labeled
501 data are primarily designed for classification and machine learning applications.
502 If the real data are unlabeled, then the evaluation of outlier detection accuracy
503 must be done based on domain knowledge or with the help of a domain expert.
504 However, this scenario is not practical for the purpose of evaluation since domain
505 knowledge is not always available. All these factors make the evaluation of the
506 proposed methodology a challenging task.

507 In this paper, we saliently illustrate the performance of our approach in
508 handling outliers using real data from the UCI Machine Learning Repository ¹.
509 Most of these data sets are labeled for classification purposes. Here, we have
510 to be aware of the fact that these class labels are not the perfect ground truth
511 in the sense that they do not correspond necessarily to potential outliers in
512 the data. Keeping these issues in mind, we have adopted a principled way to
513 produce real data for use in outlier detection.

514 In our experiments, similar to the work in (Das and Schneider, 2007), we
515 create simulated outlier objects by randomly selecting attribute values. Specif-
516 ically, in the numerical attribute space, we first normalize the attribute values
517 of each numerical attribute onto the interval $[0, 1]$ and then inject outlier
518 points whose attribute values are randomly selected from $[0, 1]$. As a result of
519 this process, all the points in our data sets have coordinates in the range $[0, 1]$
520 and are either normal points or outliers. Note that the outliers are distributed
521 at random throughout the entire space. On the other hand, to obtain outliers
522 in the categorical space, we inject novel objects in the data set in such a way
523 that, for each dimension t , the attribute value of the newly generated object is
524 randomly selected from the whole set of distinct categorical values that form

¹<http://archive.ics.uci.edu/ml/>

525 dimension t in the original data. Outliers in the mixed-attribute space are a
526 random combination of the newly generated objects in both the numerical and
527 the categorical spaces.

528 For the purpose of evaluation, we used the following standard metrics: (1)
529 Accuracy, which corresponds to the proportion of correctly partitioned objects,
530 (2) True Positive Rate (TPR), measuring the proportion of outliers that are
531 correctly identified as outliers, (3) False Positive Rate (FPR), corresponding to
532 the proportion of inliers incorrectly classified as outliers, and (4) F-measure of
533 the outliers class, corresponding to the harmonic mean between precision and
534 recall of the outlier objects class.

535 *3.2. Experiments on Mixed-Attribute Data*

536 The goal of the experiments conducted in this section are to evaluate the
537 suitability of the proposed approach in handling outliers in mixed-attribute data.
538 We compare the performance of our approach to that of ODMAD (Koufakou
539 and Georgiopoulos, 2010), the most recent approach for detecting outliers in the
540 mixed-attribute space. Note that ODMAD requires a number of parameters to
541 be set by the user. For fairness in comparison, several values were tried for the
542 parameters of ODMAD, following the suggestions in its original paper, and we
543 report results for the parameter settings that produced the best results. Note
544 that the selection of the best result here refers to the best F-measure value,
545 since this metric represents a good trade-off between TPR and FPR.

546 We considered mixed-attribute data sets taken from the UCI Machine Learn-
547 ing Repository. As mentioned in the previous subsection, to obtain data sets
548 for use in outlier detection, we generated outlier objects by randomly flipping
549 attribute values. We fixed the number of outliers injected in each set to 10%
550 of the original data set size under investigation. Fig. 5 summarizes the main
551 characteristics of the data sets used in our experiments. Note that some data
552 sets (such as Credit Approval, Automobile and Cylinder Bands) originally con-
553 tain a number of objects with missing attribute values. In our experiments, we
554 simply ignore such objects.

	#continuous attributes	#categorical attributes	#inliers	#outliers
Australian Credit Approval	6	8	690	69
German Credit	7	13	1,000	100
Credit Approval	6	9	653	65
Heart	5	8	270	27
Thoracic Surgery	3	14	470	47
Auto MPG	5	3	398	40
Automobile	15	11	159	16
Contraceptive Method Choice	2	7	1,473	147
AutoUniv (au_6)	5	7	1,000	100
Cylinder Bands	20	19	277	27

Figure 5: Mixed-attribute data sets characteristics.

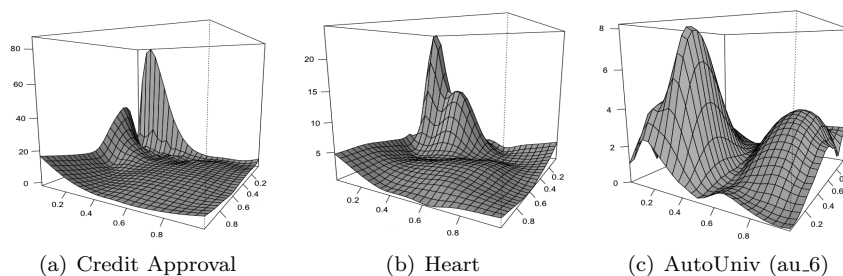


Figure 6: Estimated density curve of the outlier score vectors that correspond to three mixed-attribute data sets.

555 We used our approach to identify outliers in each of the mixed-attribute
556 data sets considered in these experiments. To this end, we set M_max to 5 and
557 then, as discussed in Section 2, we selected the optimal number of components
558 that minimize ICL-BIC. Here, the reader should be aware that the value of
559 M_max is not limited to 5 and the user can set any other value. Interestingly,
560 we found that the estimated outlier score vectors in each of the ten data sets are
561 well fitted by three bivariate beta components. For the purpose of illustration
562 and in order to not encumber the paper, we show in Fig. 6 the estimated
563 probability density function of the outlier score vectors, that corresponds to
564 Credit Approval, Heart and AutoUniv (au_6) only. Data points associated with
565 the bivariate beta component that contains the score vectors with the highest
566 values correspond to outliers. Recall that the identification of the component
567 containing the highest score values follows the procedure described in Section
568 2.3.5.

	Accuracy		TPR		FPR		F-measure	
	Proposed	ODMAD	Proposed	ODMAD	Proposed	ODMAD	Proposed	ODMAD
Australian Credit Approval	98.77%	98.94%	95.60%	94.20%	0.28%	0.57%	0.972	0.942
German Credit	98.72%	96.54%	100.00%	81.00%	1.40%	1.90%	0.934	0.810
Credit Approval	98.74%	98.60%	98.46%	92.30%	1.22%	0.76%	0.934	0.923
Heart	97.60%	93.26%	88.80%	62.96%	1.48%	3.70%	0.872	0.630
Thoracic Surgery	97.05%	94.44%	98.58%	87.94%	3.40%	3.62%	0.939	0.879
Auto MPG	90.58%	92.37%	87.50%	57.50%	9.11%	4.18%	0.625	0.575
Automobile	95.97%	94.25%	100.00%	66.60%	4.40%	3.14%	0.811	0.666
Contraceptive Method Choice	94.00%	93.20%	67.34%	62.58%	3.25%	3.73%	0.673	0.523
AutoUniv	94.54%	87.27%	88.18%	30.00%	4.82%	7.00%	0.746	0.300
Cylinder Bands	99.34%	98.68%	100.00%	92.59%	0.72%	0.72%	0.964	0.926
Average	96.53%	94.75%	92.45%	72.77%	3.01%	2.93%	0.847	0.717

Figure 7: Performance results on mixed-attribute data sets.

569 Fig. 7 compares the proposed method with ODMAD. Shaded regions in this
570 figure correspond to the best values of the four evaluation metrics considered
571 in the experiment. As can be seen from Fig. 7, our approach achieves the
572 highest true positive rates and F-measure values across all the data sets under
573 investigation and reports low false positive rates with high accuracy values. In
574 fact, the proposed method achieves, on average, an accuracy of 96.53%, TPR
575 and FPR of 92.45% and 3.01% respectively and finally an F-measure of 0.847, all
576 pointing to fairly accurate results. On the other hand, the results provided by
577 ODMAD are, on average, reasonable but less competitive than those achieved by
578 our approach. As depicted by Fig. 7, ODMAD reports, on average, an accuracy
579 of 94.75%, TPR and FPR of 72.22% and 2.93% respectively and finally an F-
580 measure of 0.717. Overall, in term of Accuracy, TP rate and F-measure, the
581 proposed method performs better than ODMAD while the FPR achieved by
582 both approaches are comparable.

583 From Fig 7, we observe that our proposed method reports an average 92.45%
584 TP rate. This means that 7.55%, on average, of outliers were misclassified as
585 inliers by our approach. This not necessarily an error, since data points have
586 coordinates in the range $[0, 1]$ and are either inliers or outliers. Outliers were
587 randomly placed throughout the entire space. In this setting, it is probable that
588 some of the outlier objects will have attribute values related to normal objects
589 in the data set under investigation. Under these circumstances, it is possible
590 that few outlier objects will have low outlier score values and consequently be

591 considered as inliers.

592 To summarize, the results presented in Fig. 7, suggest that the proposed
593 method performs well on different data sets. Furthermore, in contrast to ODMAD
594 which suffers from its dependency on several input parameters (detection thresh-
595 old, minimum support, the maximum length of itemset and the size of a window
596 of categorical and numerical scores), our approach is able to accurately iden-
597 tify outliers in an automatic fashion. Such a notable feature of our approach
598 illustrates its practical usability to effectively identify outliers in real-life appli-
599 cations. Another advantage of our approach is that it is able to handle out-
600 liers in single-type (numerical or categorical) attribute data without any feature
601 transformation, while existing methods are not able to do so. The following
602 two subsections investigate this point using real data sets characterized by only
603 numerical or categorical attributes.

604 3.3. Experiments on Numerical Data

605 The experiments described in this section aim to illustrate the capability
606 of the proposed methodology in detecting outlier objects in numerical data.
607 As discussed at the end of Section 2, when the data contains only numerical
608 attributes, we associate to each object the numerical score $\mathcal{ON}(O_i^n)$ given by
609 (1). Then, we model these scores as a mixture of univariate beta mixture.²
610 The parameters of the model $\{\lambda_m, x_m, y_m\}_{(m=1, \dots, M)}$ and the optimal number
611 of components in the mixture are estimated following the reasoning described
612 in Section 2. This process results in grouping outlier scores into several beta
613 components. Data objects associated with the beta component containing the
614 highest score values are declared outliers.

615 Fig. 8 summarizes the main characteristics of the UCI numerical data sets
616 used in the experiments. Note that, as with the experiments on mixed-attribute
617 data, we have adopted the same technique to produce outliers, that is, normal-
618 izing the attribute values between 0 and 1 and then injecting outliers in the

²To fit the beta distribution, the estimated outlier scores should be first normalized in [0,1].

	#attributes	#inliers	#outliers
Cloud	10	2,048	205
Ecoli	7	336	34
Glass Identification	9	214	21
Image Segmentation	19	2,310	231
Istanbul Stock Exchange	8	538	54
Parkinson Speech	26	1,040	104
Wine Quality - Red	12	1,599	160
Wisconsin Diagnostic Breast Cancer	30	569	57
Yacht Hydrodynamics	7	308	31
Yeast	8	1,484	148

Figure 8: Numerical data sets characteristics.

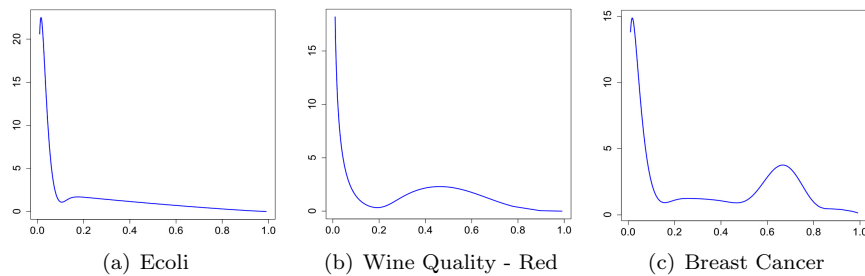


Figure 9: Estimated density curves of the numerical outlier scores that correspond to three numerical data sets.

619 data by generating objects whose attribute values are randomly selected from
620 the interval $[0,1]$. The number of outliers injected in each data set corresponds
621 to 10% of the original data set size. For each numerical data set, we estimated
622 $\mathcal{ON}(O_i^n)$ for each object and then modelled these scores as a mixture of univari-
623 ate beta distribution. To this end, we set M_{max} to 5 and selected the optimal
624 number of components that minimize ICL-BIC. We found that the number of
625 components varies from two to three beta components. For the purpose of il-
626 lustration, Fig. 9 shows the density curve of the numerical outlier scores that
627 correspond to three UCI data sets: Ecoli, Wine Quality - Red and Wisconsin
628 Diagnostic Breast Cancer. The last component in each plot depicted by Fig.
629 9 represents the highest score values. Data points associated with the scores
630 grouped in this component correspond to outliers.

631 To demonstrate the effectiveness of our approach, we compared its perfor-
632 mance to that of k NN weighed outlier algorithm (k NNW) (Angiulli and Pizzuti,
633 2005, 2002). k NNW assigns a weight to each data point based on the sum of

	Accuracy		TPR		FPR		F-measure	
	Proposed	kNNW	Proposed	kNNW	Proposed	kNNW	Proposed	kNNW
Cloud	95.47%	99.70%	94.60%	98.52%	4.40%	0.15%	0.791	0.985
Ecoli	99.18%	98.90%	93.90%	93.93%	0.29%	0.59%	0.954	0.939
Glass Identification	95.31%	93.61%	85.71%	76.19%	3.73%	4.67%	0.766	0.681
Image Segmentation	99.43%	99.74%	100.00%	98.75%	0.61%	0.14%	0.970	0.986
Istanbul stock exchange	94.39%	97.28%	81.13%	84.90%	4.29%	1.49%	0.723	0.849
Parkinson Speech	98.34%	99.39%	99.16%	96.66%	1.74%	0.33%	0.915	0.967
Wine quality - red	99.08%	98.91%	94.96%	93.71%	0.50%	0.56%	0.950	0.940
Wisconsin Diagnostic Breast Cancer	97.69%	99.18%	100.00%	98.23%	2.98%	0.50%	0.952	0.982
Yacht Hydrodynamics	97.63%	90.53%	73.33%	46.66%	0.00%	5.19%	0.846	0.467
Yeast	99.44%	97.54%	96.62%	86.48%	0.26%	1.35%	0.969	0.864
Average	97.60%	97.48%	91.94%	87.40%	1.88%	1.50%	0.884	0.866

Figure 10: Performance results on numerical data sets.

634 the distances separating that point from its k nearest neighbors in such a way
635 that outliers are characterized by high weights while inliers receive low weight
636 values. After ranking data points based on the estimated weights, the top n
637 points are identified as outliers. The implementation of this algorithm, and
638 many other outlier detection approaches, is available in the ELKI Data Min-
639 ing Framework ³ (Achtert, Kriegel, Schubert, and Zimek, 2013). Note that we
640 have chosen k NNW for its effectiveness. In fact, in our empirical investigation,
641 we have evaluated several other mainstream outlier detection algorithms, such
642 as COP (Kriegel, Kroger, Schubert, and Zimek, 2012), LDOF (Zhang, Hutter,
643 and Jin, 2009), LOCI (Papadimitriou, Kitagawa, Gibbons, and Faloutsos, 2003)
644 and LOF (Breunig et al., 2000), already implemented in ELKI. We found that
645 k NNW was the algorithm which performs well.

646 Fig. 10 illustrates the results of our approach and those of k NNW on the
647 numerical data sets considered in the experiments. Shaded regions correspond
648 to the best Accuracy, TPR, FPR and F-measure values. Recall that k NNW
649 produces a ranked list of points expecting outliers to come first. Accordingly,
650 to distinguish outliers from inliers, the user should specify the target number of
651 outliers n . In this setting, and in order to compute the value of the four evalu-
652 ation metrics used in the experiments (Accuracy, TPR, FPR and F-measure),
653 we have simply set the value of n equal to the real number of outliers in the

³<http://elki.dbs.ifi.lmu.de>

	#attributes	#inliers	#outliers
Audiology (Standardized)	69	226	23
Congressional Voting Records	16	435	43
Lymphography	18	148	15
Mushroom	22	8214	821
Primary Tumor	17	339	34
Solar Flare	10	1389	139
Soybean (Large)	35	307	31

Figure 11: Categorical data sets characteristics.

654 data set under investigation. Finally note that, as with ODMAD, we have tried
655 multiple values of k for k NNW, and we only report the best results, that is,
656 those which correspond to the highest F-measure value.

657 As can be seen from Fig. 10, our approach achieves, on average, the highest
658 Accuracy (97.60%), TPR (91.94%) and F-measure (0.884). On the other hand,
659 k NNW reports the lowest average FPR (1.50%) while our approach achieves an
660 average FPR of 1.88%. Overall, both competing algorithms show good perfor-
661 mances. A significant advantage of our approach is that it is able to automati-
662 cally discriminate outliers from inliers while with k NNW the user should specify
663 how many points should be selected as outliers.

664 3.4. Experiments on Categorical Data

665 The aim of this section is to illustrate the suitability of the proposed ap-
666 proach for handling outliers in data sets with categorical attributes only. To this
667 end, we selected a number of categorical data from the UCI Machine Learning
668 Repository. Recall that these data sets are principally labeled for classification
669 purposes. Accordingly, as discussed in Section 3.1, to produce data for use in
670 outlier detection, we inject novel data points in such a way that each attribute
671 value of each newly inserted object is randomly selected from the set of distinct
672 categorical values that initially form the corresponding attribute in the original
673 data. As with our previous experiments, the number of outliers injected in each
674 data set corresponds to 10% of the original data set size. The main characteris-
675 tics of the categorical data sets used in the experiments are summarized in Fig.
676 11.

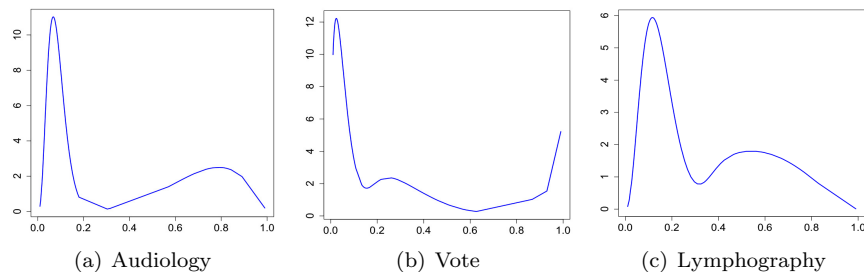


Figure 12: Estimated density curves of the categorical outlier scores that correspond to three categorical data sets.

677 To identify outliers in each of the categorical data sets considered in these
 678 experiments, we estimated first $\mathcal{OC}_{inv}(O_i^c)$ for each object. These scores are then
 679 normalized in $[0,1]$ and modelled as a mixture of univariate beta distribution.
 680 To identify the optimal number of components in the mixture, we set M_{max} to
 681 5 and selected the number of components that minimize ICL-BIC. Interestingly,
 682 as with the experiment on numeric data, we found that the optimal number of
 683 components varies from two to three. Fig. 12 illustrates the density curve of the
 684 outlier scores corresponding to three data sets: Audiology, Congressional Voting
 685 Records (Vote) and Lymphography. The last component in each plot depicted
 686 by Fig. 12 represents the highest score values. Data points associated with the
 687 scores grouped in this component correspond to outliers. The knowledgeable
 688 reader can also observe in this rendering, and also from the pervious illustration
 689 of the estimated density curves depicted by Fig. 9 and Fig. 6, the great shape
 690 flexibility of the beta distribution which leads to accurate partitioning of the
 691 outlier scores.

692 Fig. 13 compares the effectiveness of our approach to that of a recent out-
 693 lier detection approach for categorical data named Information-Theory Based
 694 Single-Pass (ITB-SP) (Wu and Wang, 2013). It has been empirically illustrated
 695 that ITB-SP is an effective approach which outperforms several existing cat-
 696 egorical outlier detection algorithms. The implementation of this algorithm
 697 has been kindly provided by its authors. As the name implies, this approach
 698 harnesses information theory concepts to estimate an outlier score for each ob-

	Accuracy		TPR		FPR		F-measure	
	Proposed	ITB-SP	Proposed	ITB-SP	Proposed	ITB-SP	Proposed	ITB-SP
Audiology (Standardized)	99.09%	100.00%	100.00%	100.00%	0.01%	0.00%	0.952	1.000
Congressional Voting Records	92.17%	92.21%	91.53%	83.07%	7.64%	5.05%	0.844	0.830
Lymphography	93.82%	90.74%	100.00%	100.00%	6.75%	10.10%	0.736	0.651
Mushroom	95.19%	98.16%	97.29%	89.90%	5.00%	1.00%	0.786	0.899
Primary Tumor	99.31%	98.62%	92.85%	92.30%	0.00%	0.75%	0.969	0.923
Solar Flare	93.05%	93.97%	87.68%	66.66%	6.40%	3.31%	0.695	0.666
Soybean (Large)	94.48%	94.48%	94.56%	88.04%	5.53%	3.58%	0.887	0.880
Average	95.30%	95.45%	94.84%	88.57%	4.48%	3.40%	0.838	0.836

Figure 13: Performance results on categorical data sets.

699 ject. Specifically, the authors in Wu and Wang (2013) propose the concept of
700 holoentropy as a new measure for outlier detection. As defined in Wu and Wang
701 (2013), holoentropy is a combination between entropy and total correlation with
702 attribute weighting, where the entropy measures the global disorder in the data
703 and the total correlation measures the attributes relationship. Based on this
704 concept, that is holoentropy, the authors formulate a function to estimate an
705 outlier score for each object in such a way that outliers are characterized by
706 high score values. The top n objects with the highest score values are declared
707 as outliers. Note that, since ITB-SP requires the number of outliers in the data
708 n to be specified by the user, and in order to compute Accuracy, TPR, FPR
709 and F-measure, we have simply set the value of n equal to the real number of
710 outliers in the data set under investigation.

711 As can be seen from Fig. 13, the average performance results for our ap-
712 proach and ITB-SP are quite similar except for the average TPR and FPR. Our
713 method reports an average 94.84% of true positives while the average TPR of
714 ITB-SP is 88.57%. This means that only 5.16%, on average, of outliers were
715 misclassified as inliers by our approach while 11.43%, on average, of outliers were
716 misclassified as inliers by ITB-SP. On the other hand, as illustrated by Fig. 13,
717 we can see that ITB-SP achieves the lowest FPR, that is 3.40%, while the pro-
718 posed method reports an average 4.48% of false positives. Overall, the results
719 illustrated in Fig. 13 suggest that both approaches display good performance.
720 Our approach has, however, the non-negligible advantage of automatically dis-
721 criminating outliers from inliers while ITB-SP requires the number of outliers in

722 the data to be specified by the user. As discussed earlier, in real applications for
723 which no prior knowledge about the data is available, it is not always possible
724 for the user to set accurately the value of this parameter.

725 **4. Conclusion**

726 In this paper, we have highlighted some limitations of existing outlier detec-
727 tion approaches for mixed-attribute data, including their dependency on user
728 parameters, such as the detection threshold and the target number of outliers to
729 be identified, which are difficult to tune and their incapability of formally dis-
730 criminating between outliers and inliers. To alleviate these problems, we have
731 proposed a principled approach that performs outlier detection in an automatic
732 fashion.

733 In our approach, we first devised two functions in order to estimate, for each
734 object, an outlier score in the numerical space and another score in the cate-
735 gorical space. Outliers in both spaces are characterized by high score values.
736 Next, we associate to each data point in the data set under investigation a two-
737 dimensional vector such that the first element of this vector corresponds to the
738 estimated outlier score in the numerical space, while the second element corre-
739 sponds to the outlier score estimated in the categorical space. Then, we model
740 these vectors as a mixture of bivariate beta. The bivariate beta component
741 that corresponds to the highest score values represents outliers. The beta dis-
742 tribution has been chosen due to its great shape flexibility which leads, in turn,
743 to accurate fitting of the estimated outlier score vectors. We have described a
744 statistical framework to illustrate how the bivariate beta mixture model can be
745 used to identify outlier objects.

746 Finally, we have devised a detailed empirical study to illustrate the suit-
747 ability of our approach in detecting outliers using several UCI data sets with
748 mixed-attributes. We have compared the performance of the proposed method
749 to that of ODMAD, the most recent approach for detecting outliers in the mixed-
750 attribute space. The results show that our approach achieves results that are,
751 in most cases, better than those of ODMAD. Moreover, we have performed

752 further experiments to demonstrate the capability of our methodology in han-
753 dling outliers in single-type attribute data without any feature transformation.
754 Tests and comparison with previous ranking approaches on several numerical
755 and categorical UCI data sets show that the proposed methodology exhibits
756 competitive results.

757 **Acknowledgements**

758 The author gratefully thanks Dr. Shu Wu for providing the implementation
759 of the ITB-SP algorithm. The author also would like to thank the reviewers for
760 their valuable comments and important suggestions. This work is supported by
761 research grants from the Natural Sciences and Engineering Research Council of
762 Canada (NSERC).

763 **References**

- 764 Achtert, E., Kriegel, H.-P., Schubert, E., & Zimek, A., 2013. Interactive data
765 mining with 3D-parallel-coordinate-trees. In Proceedings of the ACM SIG-
766 MOD international conference on Management of Data, 1009–1012.
- 767 Aggarwal, C. C., 2013. Outlier Analysis. Springer.
- 768 Alan, O., & Catal, C., 2011. Thresholds based outlier detection approach for
769 mining class outliers: An empirical case study on software measurement
770 datasets. *Expert Systems with Applications*, 38 (4), 3440–3445.
- 771 Angiulli, F., & Fassetti, F., 2014. Exploiting domain knowledge to detect out-
772 liers. *Data Mining and Knowledge Discovery*, 28 (2), 519–568.
- 773 Angiulli, F., & Pizzuti, C., 2002. Fast outlier detection in high dimensional
774 spaces. In Proceedings of the 6th European Conference on Principles of Data
775 Mining and Knowledge Discovery, 15–26.
- 776 Angiulli, F., & Pizzuti, C., 2005. Outlier mining in large high-dimensional data
777 sets. *IEEE Transactions on Knowledge and Data Engineering*, 17 (2), 203–
778 215.

- 779 Bain, L. J., & Engelhardt, M., 2000. Introduction to Probability and Mathe-
780 matical Statistics, 2nd Edition. Duxbury Press.
- 781 Bezdek, J., 1981. Pattern Recognition with Fuzzy Objective Function Algo-
782 rithms. Plenum.
- 783 Bouguessa, M., 2012. Modeling outlier score distributions. In Proceedings of
784 the 8th international conference on Advanced Data Mining and Applications,
785 713–725.
- 786 Bouguessa, M., & Wang, S., 2009. Mining projected clusters in high-dimensional
787 spaces. IEEE Transactions on Knowledge and Data Engineering, 21 (4), 507–
788 522.
- 789 Bouguessa, M., Wang, S., & Sun, H., 2006. An objective approach to cluster
790 validation. Pattern Recognition Letters, 27 (13), 1419–1430.
- 791 Bouguila, N., & Elguebaly, T., 2012. A fully Bayesian model based on reversible
792 jump MCMC and finite beta mixtures for clustering. Expert Systems with
793 Applications, 39 (5), 5946–5959.
- 794 Bouguila, N., Ziou, D., & Monga, E., 2006. Practical bayesian estimation of
795 a finite beta mixture through gibbs sampling and its applications. Statistics
796 and Computing, 16 (2), 215–225.
- 797 Boutemedjet, S., Ziou, D., & Bouguila, N., 2011. Model-based subspace clus-
798 tering of non-gaussian data. Neurocomputing, 73 (10-12), 1730–1739.
- 799 Breunig, M. M., Kriegel, H.-P., Ng, R., & Sander, J., 2000. LOF: Identifying
800 density-based local outliers. In Proceedings of the ACM SIGMOD interna-
801 tional conference on Management of Data, 93–104.
- 802 Cao, H., Si, G., Zhang, Y., & Jia, L., 2010. Enhancing effectiveness of density-
803 based outlier mining scheme with density-similarity-neighbor-based outlier
804 factor. Expert Systems with Applications, 37 (12), 8090–8101.

- 805 Das, K., & Schneider, J., 2007. Detecting anomalous records in categorical
806 datasets. In Proceedings of the 13th ACM SIGKDD international conference
807 on Knowledge Discovery and Data Mining, 220–229.
- 808 Dean, N., & Nugent, R., 2013. Clustering student skill set profiles in a unit
809 hypercube using mixtures of multivariate betas. *Advances in Data Analysis
810 and Classification*, 7 (3), 339–357.
- 811 Dempster, A., Laird, N., & Rubin, D., 1977. Maximum likelihood from incom-
812 plete data via the EM algorithm. *Journal of Royal Statistical Society, (Series
813 B)*, 39, 1–37.
- 814 Figueiredo, M. A. T., & Jain, A. K., 2002. Unsupervised learning of finite mix-
815 ture models. *IEEE Transactions on Pattern Analysis and Machine Intelli-
816 gence*, 24 (3), 381–396.
- 817 Fustes, D., Dafonte, C., Arcay, B., Manteiga, M., Smith, K., Vallenari, A., &
818 Luri, X., 2013. SOM Ensemble for unsupervised outlier analysis. Application
819 to outlier identification in the Gaia astronomical survey. *Expert Systems with
820 Applications*, 40 (5), 1530–1541.
- 821 He, Z., Xu, X., Huang, J., & Deng, S., 2005. FP-Outlier: Frequent pattern based
822 outlier detection. *Computer Science and Information System*, 2 (1), 103–118.
- 823 Huang, B., & Yang, P., 2011. Finding key knowledge attribute subspace of out-
824 liers in high-dimensional dataset. *Expert Systems with Applications*, 38 (8),
825 10147–10152.
- 826 Ji, Y., Wu, C., Liu, P., Wang, J., & Coombes, K., 2005. Applications of beta-
827 mixture models in bioinformatics. *Bioinformatics*, 21 (9), 2118–2122.
- 828 Koufakou, A., & Georgiopoulos, M., 2010. A fast outlier detection strategy for
829 distributed high-dimensional data sets with mixed attributes. *Data Mining
830 and Knowledge Discovery*, 20 (2), 259–289.

- 831 Koufakou, A., Secretan, J., & Georgiopoulos, M., 2011. Fast outlier detection in
832 large high-dimensional categorical data. *Knowledge and Information Systems*,
833 29 (3), 697–725.
- 834 Kriegel, H.-P., Kroger, P., Schubert, E., & Zimek, A., 2011. Interpreting and
835 unifying outlier scores. In *Proceedings of the 11th SIAM international con-*
836 *ference on Data Mining*, 13–24.
- 837 Kriegel, H. P., Kroger, P., Schubert, E., & Zimek, A., 2012. Outlier detection in
838 arbitrarily oriented subspaces. In *Proceedings of the 12th IEEE international*
839 *conference on Data Mining*, 379–388.
- 840 Maervoet, J., Vens, C., Berghe, G. V., Blockeel, H., & Causmaecker, P. D., 2012.
841 Outlier detection in relational data: A case study in geographical information
842 systems. *Expert Systems with Applications*, 39 (5), 4718–4728.
- 843 Otey, M. E., Ghoting, A., & Parthasarathy, S., 2006. Fast distributed outlier de-
844 tection in mixed-attribute data sets. *Data Mining and Knowledge Discovery*,
845 12 (2-3), 203–228.
- 846 Papadimitriou, S., Kitagawa, H., Gibbons, P. B., & Faloutsos, C., 2003. LOCI:
847 Fast outlier detection using the local correlation integral. In *Proceedings of*
848 *the 19th IEEE international conference on Data Engineering*, 315–326.
- 849 Penny, K., & Jolliffe, I., 2011. A comparison of multivariate outlier detection
850 methods for clinical laboratory safety data. *Journal of the Royal Statistical*
851 *Society. Series D (The Statistician)*, 50 (3), 295–308.
- 852 Tan, P.-N., Steinbach, M., & Kumar, V., 2006. *Introduction to Data Mining*.
853 Addison Wesley.
- 854 Wu, S., & Wang, S., 2013. Information-theoretic outlier detection for large-scale
855 categorical data. *IEEE Transactions on Knowledge and Data Engineering*,
856 25 (3), 589–602.

857 Yamanishi, K., Takeuchi, J.-I., Williams, G., & Milne, P., 2000. On-line un-
858 supervised outlier detection using finite mixtures with discounting learning
859 algorithms. In Proceedings of the 6th ACM SIGKDD international confer-
860 ence on Knowledge Discovery and Data Mining, 320–324.

861 Ypma, T. J., 1995. Historical development of the Newton-Raphson method.
862 SIAM Review, 37 (4), 531–551.

863 Zhang, K., Hutter, M., & Jin, H., 2009. A new local distance-based outlier
864 detection approach for scattered real-world data. In Proceedings of the 13th
865 Pacific-Asia Conference on Advances in Knowledge Discovery and Data Min-
866 ing, 813–822.

867 Zhang, K., & Jin, H., 2010. An effective pattern based outlier detection ap-
868 proach for mixed attribute data. In Proceedings of the 23rd Australasian Joint
869 Conference on Advances in Artificial Intelligence, Lecture Notes in Artificial
870 Intelligence, 6464, 122–131.