

EUR Research Information Portal

A general framework for implementing distances for categorical variables

Published in:
Pattern Recognition

Publication status and date:
Published: 01/05/2024

DOI (link to publisher):
[10.1016/j.patcog.2024.110547](https://doi.org/10.1016/j.patcog.2024.110547)

Document Version
Version created as part of publication process; publisher's layout; not normally made publicly available

Document License/Available under:
CC BY

Citation for the published version (APA):
van de Velden, M., D'Enza, A. I., Markos, A., & Cavicchia, C. (2024). A general framework for implementing distances for categorical variables. *Pattern Recognition*, 153, Article 110547. <https://doi.org/10.1016/j.patcog.2024.110547>

[Link to publication on the EUR Research Information Portal](#)

Terms and Conditions of Use

Except as permitted by the applicable copyright law, you may not reproduce or make this material available to any third party without the prior written permission from the copyright holder(s). Copyright law allows the following uses of this material without prior permission:

- you may download, save and print a copy of this material for your personal use only;
- you may share the EUR portal link to this material.

In case the material is published with an open access license (e.g. a Creative Commons (CC) license), other uses may be allowed. Please check the terms and conditions of the specific license.

Take-down policy

If you believe that this material infringes your copyright and/or any other intellectual property rights, you may request its removal by contacting us at the following email address: openaccess.library@eur.nl. Please provide us with all the relevant information, including the reasons why you believe any of your rights have been infringed. In case of a legitimate complaint, we will make the material inaccessible and/or remove it from the website.

Journal Pre-proof

A general framework for implementing distances for categorical variables

Michel van de Velden, Alfonso Iodice D'Enza, Angelos Markos,
Carlo Cavicchia



PII: S0031-3203(24)00298-X
DOI: <https://doi.org/10.1016/j.patcog.2024.110547>
Reference: PR 110547

To appear in: *Pattern Recognition*

Received date : 11 January 2023
Revised date : 17 April 2024
Accepted date : 28 April 2024

Please cite this article as: M.v.d. Velden, A.I. D'Enza, A. Markos et al., A general framework for implementing distances for categorical variables, *Pattern Recognition* (2024), doi: <https://doi.org/10.1016/j.patcog.2024.110547>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

A general framework for implementing distances for categorical variables

Michel van de Velden^a, Alfonso Iodice D'Enza^b, Angelos Markos^c, Carlo Cavicchia^d

^a*Econometric Institute, Erasmus University Rotterdam*

^b*Department of Political Sciences, University of Naples Federico II*

^c*Department of Primary Education, Democritus University of Thrace*

^d*Econometric Institute, Erasmus University Rotterdam*

A general framework for implementing distances for categorical variables

Michel van de Velden^a, Alfonso Iodice D'Enza^b, Angelos Markos^c, Carlo Cavicchia^d

^a*Econometric Institute, Erasmus University Rotterdam*

^b*Department of Political Sciences, University of Naples Federico II*

^c*Department of Primary Education, Democritus University of Thrace*

^d*Econometric Institute, Erasmus University Rotterdam*

Abstract

The degree to which objects differ from each other with respect to observations on a set of variables, plays an important role in many statistical methods. Many data analysis methods require a quantification of differences in the observed values which we can call distances. An appropriate definition of a distance depends on the nature of the data and the problem at hand. For distances between numerical variables, there exist many definitions that depend on the size of the observed differences. For categorical data, the definition of a distance is more complex as there is no straightforward quantification of the size of the observed differences. In this paper, we introduce a flexible framework for efficiently computing distances between categorical variables, supporting existing and new formulations tailored to specific contexts. In supervised classification, it enhances performance by integrating relationships between response and predictor variables. This framework allows measuring differences among objects across diverse data types and domains.

Keywords: Categorical data, Distance, Cluster analysis, Classification, K-NN

1 **1. Introduction**

2 In many statistical methods, the quantification of dissimilarity, that is, the
3 degree to which objects differ from each other, plays an important role [1]. We can
4 refer to such dissimilarity quantification as a distance. Classification methods such
5 as *K*-Nearest Neighbors [KNN, 2], but also clustering methods such as *K*-means,
6 Partitioning Around Medoids (PAM) and hierarchical linkage methods [3], and
7 data visualization methods such as MultiDimensional Scaling [MDS, 4] and biplots
8 [5], require a definition of distance between objects. The way to select a definition
9 of distance depends on the nature of the data and problem at hand.

10 The distance measures for numerical data are usually based on the magnitude
11 of the observed differences in values. For categorical data, however, the situation
12 is more complex as we do not directly observe, and hence cannot directly quantify,
13 sizes of differences. We can only establish whether there is a difference or not.
14 For distance calculations in multivariate contexts, two cases can be distinguished.
15 First, the distances are calculated for each variable independently and then added.
16 Second, the association between the variables is taken into account when calcu-
17 lating the distances. For numerical variables, several well-known distances, e.g.,
18 Euclidean or Manhattan distances, implicitly assume independence between the
19 variables and require that the measurement scales must be commensurable. For
20 categorical variables, there are also several measures that take the sum of distances
21 per variable when considering a multivariate distance. For example, in simple
22 matching, distance between two observations is defined as the number of times
23 that the categories of corresponding variables do not match.

24 For numerical variables, the association between variables can be accounted for
25 using the Mahalanobis distance, where (sample) covariances are used to weigh ob-
26 served differences to account for correlation between the variables. For categorical
27 data, so-called association-based distances exist. In such distances, the association
28 between categorical variables is used to quantify the differences between observa-
29 tions. The question of how to account for associations in a categorical setting is
30 not trivial. Several recent proposals for distances between categorical variables are
31 indeed association-based distances [see, e.g., 6, 7, 8, 9, 10, 11].

32 In this paper, we propose a general framework for implementing categorical dis-
33 tances. By incorporating existing distances in our framework, it becomes possible
34 to assess the differences and similarities between them. Currently, such comparison
35 is not trivial due to the wide variety of notation and research fields (and hence
36 objectives) in which methods have been proposed. In addition, our framework
37 makes it possible to construct and define new and highly customizable distances.
38 For example, in a supervised classification context, the framework can be used to
39 define a distance that takes into account association with the classes of the response
40 variable. As our framework is not method- or application-specific, it can be used
41 to calculate distance matrices for any method or application that requires distance
42 calculations. For example, MDS, cluster analysis, or, in a supervised context, KNN.
43 The framework offers flexibility and transparency by defining the categorical vari-
44 able distances as a function of category dissimilarities. The category dissimilarities
45 can either be data-driven or theory-based. Furthermore, our framework allows
46 for fast and efficient distance calculations, potentially representing a substantial

47 enhancement compared to current implementations. In particular, we show that for
48 the distance for categorical variables proposed in [7], our implementation is much
49 faster than existing implementations.

50 An important issue with regard to the definition of distance measures is the
51 validation of the measures. That is, how does one know that the chosen measure
52 is appropriate? Although the new framework does not provide an answer to this
53 question, having a general formulation simplifies both theoretical and empirical
54 comparisons between different choices. We illustrate our method by applying
55 distance-based data analysis methods to several well-known categorical data sets
56 using a selection of known and new, association-based, categorical distances. We
57 implemented functions to perform all categorical distance calculations using our
58 general framework in the R package `catdist`, which is available on GitHub¹.

59 This paper is organized as follows. After introducing some notation in Section
60 2, we describe our general framework in Section 3. In Section 4, we introduce
61 several common distance measures for categorical variables and show how they
62 can be incorporated. Categorical distances based on co-occurrences are introduced
63 in Section 5, with particular attention to the distance measure proposed by [7]. In
64 Section 6, we show how supervised distances can be constructed and implemented
65 using our framework. Tuning of distance definitions is described in Section 7, after
66 which we illustrate our methodology using several data sets in Section 8. Section 9
67 concludes the paper.

¹https://github.com/alfonsoIodiceDE/catdist_package

68 **2. Notation**

Suppose that we have n observations on Q categorical variables and let the number of categories for the $j \in 1 \dots Q$ -th variable be q_j . We can then code the categorical data by using indicator matrices. That is, for each categorical variable $j \in 1 \dots Q$, we create an $n \times q_j$ binary matrix \mathbf{Z}_j , where the n rows correspond to observations and the q_j columns to categories. The observed category is indicated by a one, and all other categories are assigned zeros. Furthermore, for each observation of a categorical variable, exactly one category is observed, and we only include categories that have been observed at least once in the data set. Therefore, each column of \mathbf{Z}_j contains at least one element equal to one and $\mathbf{Z}_j \mathbf{1}_{q_j} = \mathbf{1}_n$, where, generically, $\mathbf{1}_i$ denotes an i -dimensional vector of ones. That is, the sum over the columns is 1. Using these indicator matrices, we can code data on Q categorical variables into a so-called super-indicator matrix by collecting all indicator matrices next to each other. That is, $\mathbf{Z} = (\mathbf{Z}_1 \dots \mathbf{Z}_Q)$. Furthermore, define

$$\mathbf{P} = \frac{1}{n} \mathbf{Z}' \mathbf{Z}, \quad (1)$$

69 and $\mathbf{P}_d = \frac{1}{n} (\mathbf{Z}' \mathbf{Z}) \odot \mathbf{I}_{Q^*}$, where \odot indicates the Hadamard product, that is, element-
70 wise multiplication, and $Q^* = \sum_{j=1}^Q q_j$. Note that \mathbf{P}_d is a diagonal matrix with as
71 its diagonal elements the observed relative frequencies (within each variable) for
72 the categories. Moreover, let $\mathbf{p} = \mathbf{P}_d \mathbf{1}_{Q^*}$ denote the vector of observed relative
73 frequencies, and $\mathbf{p}^- = \mathbf{P}_d^{-1} \mathbf{1}_{Q^*}$ is the vector of inverse observed relative frequencies.

Note that the ij -th off-diagonal block of \mathbf{P} gives the relative frequencies of

co-occurrences for the categories of variables i and j . They can be seen as (empirical) joint probability distributions for variables i and j . For the calculation of association-based distances in Section 5, we also define

$$\mathbf{R} = \mathbf{P}_d^{-1} (\mathbf{P} - \mathbf{P}_d). \quad (2)$$

74 The rows of the ij -th off-diagonal block of \mathbf{R} give, for the categories of the i -th
75 variable, the distributions over the categories of the j -th variable. These can be
76 interpreted as (empirical) conditional distributions.

77 3. Categorical distance calculations based on category dissimilarities

78 For a categorical variable, it is not obvious how to quantify differences between
79 different categories. For example, suppose that we observe three individuals, one
80 from the Netherlands, one from Italy, and one from Greece. Geographically, and
81 perhaps also culturally, Italy and Greece are more similar than the Netherlands.
82 How to take such differences into account is, however, not trivial. In our framework,
83 we do so by defining *category dissimilarities*.

84 **Definition 1.** A matrix Δ_j is the category dissimilarity matrix for variable j . The
85 elements of this matrix, δ_{ab} , where a and b indicate two categories of variable j ,
86 quantify the dissimilarities between the categories a and b of the j -th variable.

87 We can impose conditions on the dissimilarity matrix that are consistent with
88 typical distance definitions. That is, 1) the dissimilarity of a category from itself is
89 zero ($\delta_{aa} = 0$, for all categories). 2) Dissimilarities are symmetric ($\delta_{ab} = \delta_{ba}$, for all
90 pairs of categories). 3) Dissimilarities satisfy the triangle inequality. That is, if a, b

91 and c denote different categories for a variable j , then, for all categories a, b and c ,
 92 $\delta_{ac} \leq \delta_{ab} + \delta_{bc}$. If all three of these conditions are satisfied, the dissimilarities can
 93 be considered as metric distances between categories. If they are non-negative and
 94 satisfy only the first two conditions, they can be interpreted as non-metric distances
 95 between categories. However, we refer to them as category dissimilarities and
 96 reserve the term “distance” for the distances between observations.

97 If we have Q categorical variables, each with a category dissimilarity matrix
 98 Δ_j , we can construct a $Q^* \times Q^*$, block diagonal matrix $\mathbf{\Delta}$, with separate category
 99 dissimilarity matrices as diagonal blocks.

100 The category dissimilarity matrices can be used to calculate a between obser-
 101 vations distance matrix as follows. First, consider the n by q_j indicator matrix
 102 \mathbf{Z}_j corresponding to the j -th categorical variable. Furthermore, we have the cor-
 103 responding category dissimilarity matrix Δ_j . We can formulate the following
 104 theorems:

105 **Theorem 1.** *The distances between the observations for the categorical variable*
 106 *j are $\mathbf{D}_j = \mathbf{Z}_j \Delta_j \mathbf{Z}_j'$.*

107 **PROOF.** The matrix multiplication of the row i of \mathbf{Z}_j with Δ_j selects the row of
 108 Δ_j corresponding to the category chosen by the individual i . Similarly, matrix
 109 multiplication of this row by the i' -th column of \mathbf{Z}_j' (i.e., the i' -th observation)
 110 selects the element corresponding to the category chosen by the individual i' .
 111 Hence, the (i, i') -th element of \mathbf{D}_j is the dissimilarity between the categories
 112 chosen by individuals i and i' . ■

113 **Theorem 2.** *If we define the distance between observations on Q categorical*
 114 *variables as the sum of Q distances for each categorical variable, the $n \times n$ distance*
 115 *matrix can be calculated as*

$$\mathbf{D} = \mathbf{Z} \mathbf{\Delta} \mathbf{Z}' . \quad (3)$$

PROOF.

$$\mathbf{D} = (\mathbf{Z}_1 \dots \mathbf{Z}_Q) \begin{pmatrix} \Delta_1 & & \\ & \dots & \\ & & \Delta_Q \end{pmatrix} \begin{pmatrix} \mathbf{Z}'_1 \\ \vdots \\ \mathbf{Z}'_Q \end{pmatrix} = \sum_{j=1}^Q \mathbf{Z}_j \Delta_j \mathbf{Z}'_j = \sum_{j=1}^Q \mathbf{D}_j$$

116

■

117 From (3), it follows that distances between observations of categorical variables
 118 depend on the choices of the category dissimilarity matrices Δ_j . This allows for
 119 great flexibility in defining a suitable distance measure for a set of categorical
 120 variables. In the next section, we briefly review some choices for Δ , and we show
 121 how they relate to existing distances.

122 Note that associations between categorical variables are not explicitly incorpo-
 123 rated in this formulation. That is, the differences between the categories observed
 124 for one variable are not related to the differences in the categories observed for
 125 other variables. There are, however, ways to account for such observations. For
 126 example, rather than creating an indicator matrix for each categorical variable,
 127 one could construct an indicator matrix for all possible combinations (or subsets
 128 thereof) of observations. That is, one can create, for each (or a subset of) combina-
 129 tion of categories one indicator matrix. The number of columns of such a matrix
 130 is therefore $\prod_{j=1}^Q q_j$ and only one category dissimilarity matrix is needed where
 131 each category is a combination of the categories for all Q variables. However, with
 132 several categorical variables, the total number of combinations and hence the num-
 133 ber of categories of the final indicator matrix quickly becomes large. Furthermore,
 134 finding an appropriate category dissimilarity matrix for the combinations is not a
 135 trivial task. An alternative way to account for associations between the categorical

136 variables is to use them in the construction of the category dissimilarity matrices.
 137 That is, by defining the dissimilarities between the categories of a variable in Δ_j ,
 138 based on the associations with other variables. In Section 5, we give examples of
 139 such category dissimilarity measures.

140 3.1. Distances between sets

141 Suppose that we have two separate sets of observations on the same Q categori-
 142 cal variables. Data for these two sets can be collected in the $n_1 \times Q^*$ and $n_2 \times Q^*$
 143 super indicator matrices $\mathbf{Z}^{(1)}$ and $\mathbf{Z}^{(2)}$. Then, for a known category dissimilarity
 144 matrix $\mathbf{\Delta}$, it is easily verified that the distances between the observations for the
 145 two sets can be calculated as $\mathbf{D}^{(12)} = \mathbf{Z}^{(1)} \mathbf{\Delta} \mathbf{Z}^{(2)T}$.

146 Note that the matrix $\mathbf{D}^{(12)}$ is of order $n_1 \times n_2$. Calculating distances between sets
 147 can be useful when considering distance-based classification problems. In KNN,
 148 for example, distances between “new” (unlabeled) observations and observations
 149 in a labeled data set are required. The KNN predictions are based on (usually by
 150 majority vote) the labels of the K nearest neighbors in the training set. Similarly,
 151 in PAM – a popular clustering algorithm similar to K-means, where instead of
 152 considering within-cluster variation around the mean, variation around an actual
 153 observation, the medoid, is considered – computing the distance between sets can
 154 be of interest. If the medoids are collected in $\mathbf{Z}^{(1)}$ and “new” data points in $\mathbf{Z}^{(2)}$,
 155 we can assign the new points to existing clusters considering the distances $\mathbf{D}^{(12)}$
 156 and selecting the smallest distances.

157 **4. Independent category dissimilarity matrices**

158 We first consider several definitions of dissimilarity for categorical data that do
 159 not take into account the association between variables. Hence, in a multivariate
 160 context, distances are calculated as the sum of distances per variable, and for
 161 each variable, the category dissimilarities are independent of the observations on
 162 other variables. However, category dissimilarities may depend on the observed
 163 frequencies for a variable. We do not aim to be complete with respect to the
 164 different definitions of dissimilarity. Instead, we select definitions from [12]
 165 [which contain several definitions also reviewed in 13], and transform them into
 166 dissimilarities. We show how these dissimilarities can be incorporated into our
 167 framework by defining the appropriate category dissimilarity matrices $\mathbf{\Delta}$. For a
 168 more detailed description, as well as study on the relative performances of these
 169 dissimilarity definitions in a cluster analysis setting, see [12].

170 *Simple Matching (SM)*. The idea of this measure is that the distance between
 171 observations is 1 if the categories do not match and 0 if they do. Consequently, all
 172 different between category dissimilarities are 1. For the j -th categorical variable
 173 with q_j categories, we define $\mathbf{\Delta}_{M_j} = \mathbf{1}_{q_j} \mathbf{1}'_{q_j} - \mathbf{I}_{q_j}$.

174 That is, the distance between each category is exactly 1. If we have Q categori-
 175 cal variables and want to calculate a simple matching for all variables, we simply
 176 collect all $\mathbf{\Delta}_{M_j}$ in a block diagonal matrix $\mathbf{\Delta}_M = \mathbf{K}_b - \mathbf{I}_{Q^*}$, where \mathbf{K}_b is a $Q^* \times Q^*$
 177 block diagonal matrix with, for $j = 1 \dots Q$, $q_j \times q_j$ matrices (\mathbf{K}_{b_j}) of ones as its
 178 diagonal blocks.

179 *Eskin*. For Eskin distance [14], category dissimilarities depend on the number

180 of categories. Dissimilarities for variables with more categories are smaller than
 181 dissimilarities for variables with fewer categories. In particular, the dissimilarity
 182 between different categories for a variable with q_j categories is $\frac{2}{q_j}$. Therefore, for
 183 the j -th categorical variable with q_j categories, the category dissimilarity matrix
 184 is defined as $\mathbf{\Delta}_{E_j} = \frac{2}{q_j} (\mathbf{1}_{q_j} \mathbf{1}'_{q_j} - \mathbf{I}_{q_j})$. Collecting all Q dissimilarity matrices in a
 185 block diagonal matrix produces the Eskin category dissimilarity matrix $\mathbf{\Delta}_E$. If all
 186 variables have the same number of categories, Eskin merely re-scales the simple
 187 matching dissimilarity.

188 *Lin.* [15] proposed an information-theoretic measure that gives more weight to
 189 matches on frequent values and lower weight to mismatches on infrequent values.
 190 We implement Lin's proposal as follows: define $\mathbf{P}_r = \mathbf{p} \mathbf{1}'_{Q^*}$ and $\mathbf{P}_c = \mathbf{1}_{Q^*} \mathbf{p}'$. Further-
 191 more, let $\tilde{\mathbf{P}} = \mathbf{P}_r + \mathbf{P}_c - \mathbf{P}_d$. Then, the category dissimilarity matrix can be defined
 192 as $\mathbf{\Delta}_{Lin} = \left[\log(\mathbf{P}_r) + \log(\mathbf{P}_c) - 2\log(\tilde{\mathbf{P}}) \right] \oslash 2\log(\mathbf{P}_r + \mathbf{P}_c)$, where \oslash indicates the
 193 Hadamard division (i.e., element-wise) and $\log(\cdot)$ takes the logarithms of the ele-
 194 ments of the parenthesized object and collects them in an object of the same size.
 195 Note that Lin's dissimilarity for a category with itself is zero. Furthermore, in our
 196 implementation, for each variable, the dissimilarity between different categories,
 197 say categories a and b , is $[\log(p_a) + \log(p_b) - 2\log(p_a + p_b)] / 2\log(p_a + p_b)$,
 198 where p_a and p_b are, respectively, the relative frequencies of categories a and
 199 b .

200 *Inverse Occurrence Frequency (IOF) and Occurrence Frequency (OF).* For
 201 inverse occurrence frequency [IOF, 13], a higher dissimilarity is assigned when
 202 categories are more frequently observed. In particular, the category dissimilarity

203 matrix is defined as $\Delta_{IOF} = [\log(np)] [\log(np)]' - [\log(np)] [\log(np)]' \odot \mathbf{I}_{Q^*}$. It
 204 is worth observing that, for each variable, IOF dissimilarity for a category with
 205 itself is zero, and the dissimilarity between two different categories, say a and b ,
 206 corresponds to $\log(np_a) \log(np_b)$. The IOF measure is related to the concept of
 207 inverse document frequency (TF-IDF) from information retrieval, where it is used
 208 to account for document relevance for a given term [16]. In other words, since a
 209 rare term contributes more information than a more frequent term, the IOF measure
 210 accounts for how rare the term is, and a lower IOF dissimilarity corresponds to
 211 a rarer term. Log frequency is used to reduce the impact of terms of very high
 212 frequencies. For occurrence frequency (OF) dissimilarity, dissimilarities are higher
 213 if the categories are observed less frequently. The category dissimilarity matrix
 214 is defined as $\Delta_{OF} = [\log(\mathbf{p}^-)] [\log(\mathbf{p}^-)]' - [\log(\mathbf{p}^-)] [\log(\mathbf{p}^-)]' \odot \mathbf{I}_{Q^*}$. Therefore,
 215 OF dissimilarity for a category with itself is zero, and the dissimilarity between
 216 two different categories a and b is $\log(p_a) \log(p_b)$.

217 *Goodall dissimilarities.* In [13], four variations of Goodall's similarity are
 218 considered. These are based on Goodall's original proposal [17]. After transform-
 219 ing similarities into dissimilarities, where dissimilarity is $1 - \text{similarity}$, the four
 220 measures have in common that dissimilarities between different categories are, as
 221 is the case with simple matching, always equal to one. However, the dissimilarity
 222 of a category with respect to the same category depends on the observed propor-
 223 tions of the categories. Below we provide the category dissimilarity matrices for
 224 Goodall 3 and Goodall 4. For Goodall 1 and 2, we can also construct such matrices.
 225 However, these definitions require conditional sums of proportions. In particular,

226 for Goodall 1, the dissimilarity for category a with itself, is defined as the sum of
 227 squared observed proportions that are smaller or equal to the observed proportion
 228 of category a . For Goodall 2, it is the sum of squared observed proportions that
 229 are larger or equal to the observed proportion of category a . The Goodall 3 and 4
 230 measures do not require the calculation of a (conditional) sum and have the squared
 231 proportion and one minus the squared proportion of a category, respectively, on
 232 the diagonal blocks of $\mathbf{\Delta}$. That is, $\mathbf{\Delta}_{G_3} = \mathbf{K}_b - \mathbf{I}_{Q^*} + \mathbf{P}_d^2$, and $\mathbf{\Delta}_{G_4} = \mathbf{K}_b - \mathbf{P}_d^2$. In
 233 these definitions, dissimilarity of a category with the same category is not zero.
 234 Consequently, the resulting “distances” do not satisfy the typical requirements of
 235 a distance. Note that for the Goodall 1 and 3 measures, a higher dissimilarity is
 236 assigned when the matching categories are frequent, whereas for the Goodall 2
 237 and 4 measures a higher dissimilarity is assigned when the matching categories are
 238 infrequent.

239 *Variable Entropy (VE) and Variable Mutability dissimilarities (VM).* [12]
 240 proposed two variability-based dissimilarity measures that are related to Goodall
 241 1 and 2, respectively. These dissimilarities are equal to one if the categories do
 242 not match, while, if they do match, VE uses the entropy and VM uses the Gini
 243 coefficient to quantify “dissimilarity”. In particular, for the j -th categorical variable
 244 with q_j categories, the category dissimilarity matrices are defined as $\mathbf{\Delta}_{VE_j} = \mathbf{K}_{b_j} +$
 245 $\left(\frac{1}{\log q_j} \sum_{l=1}^{q_j} p_l \log p_l \right) \mathbf{I}_{q_j}$ and $\mathbf{\Delta}_{VM_j} = \mathbf{K}_{b_j} - \left[\frac{q_j}{q_j-1} \left(1 - \sum_{l=1}^{q_j} p_l^2 \right) \right] \mathbf{I}_{q_j}$, respectively.
 246 Collecting all Q dissimilarity matrices in a block diagonal matrix returns the VE
 247 and VM category dissimilarity matrices $\mathbf{\Delta}_{VE}$ and $\mathbf{\Delta}_{VM}$.

248 *4.1. Ordered categories*

249 If categories are ordered, the order can be reflected in the dissimilarities. A
 250 simple choice would be to define the dissimilarities as the difference in category
 251 numbers. That is, the dissimilarity between categories a and b is simply $b - a$. If
 252 the data are rank order data or rating (e.g., Likert) scale data, this definition would
 253 imply treating the data as interval data. However, implementation of alternative,
 254 custom, definitions of ordered between-category distances is also straightforward.
 255 For example, if the categories correspond to bins on a numerical scale (e.g., age
 256 or income groups), differences between the midpoints of the bins can be used to
 257 define dissimilarities that better reflect the underlying values. More generally, let Δ_o^i
 258 denote the i -th diagonal block of Δ_o^* , and δ_{ab}^i its ab -th element. Then, dissimilarities
 259 between ordered categories can be imposed by letting $\delta_{ab}^i \leq \delta_{ab^*}^i$, for $a, b \in 1 \dots i_{q_i}$,
 260 $b^* \neq b$ and $b^* > a$. Note that this definition does not guarantee that the triangle
 261 inequality holds. That is, it may be the case that the direct distances between
 262 two categories are larger than the indirect distances between those categories. To
 263 enforce the triangle inequality to hold, one could impose appropriate additional
 264 constraints on the elements of the category dissimilarity matrix.

265 **5. Association-based category dissimilarity matrices**

266 The general idea of association-based measures, similar to the case of the
 267 Mahalanobis distance for numerical variables, differences that are in line with
 268 the association between variables are less informative (i.e., should correspond to
 269 smaller dissimilarity values) than differences that are not in line with the general

270 association. How to exactly implement this idea depends on the calculation of the
 271 association between categorical variables, and how to incorporate this association
 272 in the category dissimilarities. Here, we present a general form to calculate and
 273 collect association-based dissimilarities that can be directly implemented in our
 274 general framework in Section 3. We then present some specific variants and link
 275 them to recent proposals.

276 Fundamental in the calculation of association-based distances is the matrix
 277 of proportions of co-occurrences \mathbf{P} and the corresponding profile matrix \mathbf{R} as
 278 defined in Section 2. In particular, recall that the off-diagonal blocks of \mathbf{P} and \mathbf{R}
 279 can be interpreted as (empirical) joint and conditional probability distributions,
 280 respectively. By considering different ways to quantify the dissimilarities between
 281 the conditional distributions (i.e., the rows of the off-diagonal blocks of \mathbf{R}) we can
 282 construct different category dissimilarity matrices $\mathbf{\Delta}$ that, by applying (3) can be
 283 used to obtain the between-observation distances.

Let \mathbf{R}^{ij} denote the ij -th off-diagonal block of \mathbf{R} , and let \mathbf{r}_a^{ij} denote its a -th row. Note that the elements of each row of \mathbf{R}^{ij} add up to 1. Hence, these elements can be seen as (empirical) conditional probabilities. We define the dissimilarities between categories for all pairs of categories of variable i (for $i = 1 \dots Q$) based on the association with variable j (with $j \neq i$), as

$$\delta^{ij}(a, b) = \Phi^{ij}(\mathbf{r}_a^{ij}, \mathbf{r}_b^{ij}), \quad (4)$$

where, generically, a and b indicate categories of variable i and Φ^{ij} is the dis-

similarity function that quantifies the differences between profiles based on the association between variables i and j . The overall between-category dissimilarities for all pairs of categories of variable i can be defined as

$$\delta^i(a, b) = \sum_{j \neq i}^Q w_{ij} \delta^{ij}(a, b) = \sum_{j \neq i}^Q w_{ij} \Phi^{ij}(\mathbf{r}_a^{ij}, \mathbf{r}_b^{ij}). \quad (5)$$

284 The weights w_{ij} in (5) allow flexibility with respect to the importance of different
 285 variables in the calculation of association-based category dissimilarities, as defined
 286 by Φ^{ij} . By collecting, for each variable i , the elements $\delta^i(a, b)$ in a category
 287 dissimilarity matrix $\mathbf{\Delta}_\Phi^i$, and organizing them on the diagonal of a block diagonal
 288 matrix, we obtain a dissimilarity matrix $\mathbf{\Delta}_\Phi$, which can be used to calculate the
 289 distances between the observations using (3). Thus, 5 provides a very general
 290 way to define category dissimilarities using pair-specific weights and dissimilarity
 291 functions.

292 For the association-based dissimilarity functions Φ^{ij} , any function that quan-
 293 tifies the difference between two distributions can be used. A brief overview of
 294 46 different functions and their implementation in the R package `philentropy`
 295 is described in [18]. Concerning the choice of weights w_{ij} , we distinguish two
 296 options. In the first, all weights are equal. Usually, 1 or $1/(Q-1)$, so that sums or
 297 averages are obtained. Alternatively, different weights can be used for different
 298 pairs. These weights can either be selected using expert knowledge (e.g., based
 299 on the experience and preferences of the researcher) or by using a data-driven
 300 approach. For example, one could set certain weights to zero and others to some

301 constant based on some predetermined data dependent criterion (e.g., a measure of
 302 association like Cramér's V). Pairs with non-zero weights can then be referred to
 303 as "context" variables. Approaches using such context-based dissimilarities are
 304 described in [8, 9]. If an objective measure of overall fit of a solution is available,
 305 one could consider w_{ij} and Φ^{ij} as tuning parameters, and search combinations
 306 of these parameters to make a choice. In Section 7 we shall further explore the
 307 tuning of w_{ij} and Φ^{ij} . In the next subsections, we describe some specific choices
 308 of dissimilarity functions. In particular, we provide definitions for category dissim-
 309 ilarities between categories a and b of variable i , based on the association between
 310 variables i and j . That is, we present specific choices of Φ^{ij} in (4). Inserting
 311 these definitions into (5) results in a dissimilarity matrix Δ_{Φ} that can be used to
 312 calculate the between-observation distances. For ease of notation, we now drop the
 313 superscripts ij .

Total variation distance between profiles (TVD). The total variation distance
 between two discrete probability distributions can be defined as $1/2$ times the
 L_1 norm between the distributions. We can implement this in our framework by
 defining the category dissimilarity function Φ as

$$\Phi(\mathbf{r}_a, \mathbf{r}_b) = \frac{1}{2} \sum_{l=1}^{q_j} |r_{al} - r_{bl}|, \quad (6)$$

314 where r_{al} and r_{bl} denote the l -th element of \mathbf{r}_a and \mathbf{r}_b , respectively.

315 *Ahmad and Dey's categorical variable distance.* [7] argue that the dissimilarity
 316 between categories should be computed as a function of their distribution in

317 the overall data set and in co-occurrence with other categories, rather than in
 318 isolation. The idea is to take into account co-occurrences of categories when
 319 constructing distances. The way they do this, is by considering all combinations of
 320 categories of one variable, and selecting the partitioning (that is, a combination of
 321 categories) for which the sum of proportions in the two complementary partitions
 322 for the two categories is maximal. Following [7], we can define the dissimilarity
 323 between categories a and b of variable i , with respect to the distribution over the
 324 categories of variable j , as $\delta(a, b) = \max_{\omega_j} (P(\omega_j|a) + P(\bar{\omega}_j|b) - 1)$, where ω_j
 325 and its complement $\bar{\omega}_j$ define a binary partition with respect to the categories of
 326 variable j , and $P(\omega_j|a)$ denotes the proportion of observations with the category
 327 a of variable i , corresponding to the set of categories of variable j as defined
 328 by ω_j . Note that the term -1 is only introduced to fix the upper limit of the
 329 dissimilarities at 1. The number of binary partitions for variable j , excluding the
 330 partitions containing all or no categories, equals $2^{q_j} - 2$, where q_j gives number
 331 of categories of variable j , and hence this number grows exponentially when the
 332 number of categories for a variable increases. [7] propose an algorithm to calculate
 333 their distances. The order of their algorithm is $O(Q^*2n + Q^*2\bar{q}^3)$, where Q^* gives
 334 the total number of categories, n is the number of observations and \bar{q} denotes the
 335 average number of categories per variable.

336 The Ahmad and Dey's categorical variable distance [7] is equivalent to TVD,
 337 and computing the latter is much more efficient. However, as this relationship is
 338 not trivial and appears to be unknown, we present this below by formulating the
 339 following theorem.

340 **Theorem 3.** *Ahmad and Dey's distance for categorical variables [7] is equivalent*
 341 *to the total variation distance between profiles.*

342 **PROOF.** When going from $\delta(a, b)$ to $\delta(b, a)$, the optimal partition ω_j , that is, the
 343 combination of categories that maximizes the sum, is simply flipped (i.e., the
 344 complement is taken), hence Ahmad and Dey's distance is symmetric: $\delta(a, b) =$
 345 $\max_{\omega} (P(\omega|a) + P(\bar{\omega}|b) - 1) = \max_{\omega} (P(\omega|b) + P(\bar{\omega}|a) - 1) = \delta(b, a)$, where, for
 346 convenience, we dropped the subscripts j for the ω 's. As $P(\bar{\omega}|\cdot) = 1 - P(\omega|\cdot)$,
 347 we have $\delta(a, b) = \max_{\omega} (P(\omega|a) - P(\omega|b)) = \max_{\omega} (P(\omega|b) - P(\omega|a))$. This
 348 shows that Ahmad and Dey's distance is equal to finding the maximum differ-
 349 ence between all combinations of observed proportions. This implies that we
 350 can express this distance as the supremum norm of a vector of differences be-
 351 tween probabilities. The total variation distance can also be defined as the largest
 352 difference between probabilities from two probability distributions that can be
 353 assigned to the same event. Let \mathbf{K}_{q_j} be a design matrix that defines all, except the
 354 empty and complete, binary partitions for the q_j categories of variable j . Hence,
 355 \mathbf{K}_{q_j} is a $q_j \times q_j^*$ matrix of zeros and ones, where $q_j^* = \sum_{l=1}^{q_j-1} \binom{q_j}{l} = 2^{q_j} - 2$. Then
 356 $\delta(a, b) = \|\mathbf{K}'_{q_j}(\mathbf{r}_a - \mathbf{r}_b)\|_{\infty} = \|\mathbf{K}'_{q_j} \mathbf{d}_{ab}\|_{\infty}$, where $\mathbf{d}_{ab} = (\mathbf{r}_a - \mathbf{r}_b)$ and $\|\mathbf{x}\|_{\infty}$ denotes
 357 the supremum norm of vector \mathbf{x} , that is, the maximum element of \mathbf{x} in absolute
 358 value. The number of columns of \mathbf{K}_{q_j} and hence the size of the vector from which
 359 we need to take the norm, grows exponentially with the number of categories. A
 360 more efficient way to calculate the distances between categories a and b can be
 361 obtained using

$$\|\mathbf{K}'_{q_j} \mathbf{d}_{ab}\|_{\infty} = \frac{1}{2} \|\mathbf{d}_{ab}\|_1, \quad (7)$$

362 where $\|\mathbf{x}\|_1$ denotes the L_1 norm of vector \mathbf{x} , that is $\|\mathbf{x}\|_1 = \sum_i |x_i|$. To see that
 363 (7) holds, note that the maximum, in absolute value, for the combinations of el-
 364 ements of \mathbf{d}_{ab} is obtained by selecting the combination consisting of elements
 365 that have the same sign. As $\mathbf{d}_{ab} \mathbf{1}_{q_j} = (\mathbf{r}_a - \mathbf{r}_b) \mathbf{1}_{q_j} = 0$, the sum of all posi-
 366 tive elements equals the sum of all negative values. Therefore, $\|\mathbf{K}'_{q_j} \mathbf{d}_{ab}\|_{\infty} =$
 367 $\sum_{l: d_l > 0} \|d_l\| = \sum_{l: d_l < 0} \|d_l\|$, where d_l denotes the l -th element of \mathbf{d}_{ab} . Finally, as
 368 $\|\mathbf{d}_{ab}\|_1 = \sum_l \|d_l\| = \sum_{l: d_l > 0} \|d_l\| + \sum_{l: d_l < 0} \|d_l\|$, the equivalence in (7) immediately
 369 follows and we can express Ahmad and Dey's distance between categories a and
 370 b with respect to the categories of variable j , as $\delta(a, b) = \frac{1}{2} \|\mathbf{d}_{ab}\|_1 = \frac{1}{2} \|\mathbf{d}_{ab}\|_1 =$
 371 $\frac{1}{2} \sum_{l=1}^{q_j} |r_{al} - r_{bl}|$.

372 ■

373 *Kullback-Leibler divergence between profiles* (KL). Kullback-Leibler diver-
 374 gence [19] is an entropy-based measure of dissimilarity between probability dis-
 375 tributions. [6] define category dissimilarities for the categories of variable i by
 376 taking the sum of KL-divergences between the (empirical) conditional proba-
 377 bility distributions for all other variables. Using similar notation as before, we
 378 can implement this divergence by setting all weights w_{ij} equal to one and by
 379 defining $\Phi(\mathbf{r}_a, \mathbf{r}_b) = \sum_{l=1}^{q_j} \left[r_{al} \log\left(\frac{r_{al}}{r_{bl}}\right) + r_{bl} \log\left(\frac{r_{bl}}{r_{al}}\right) \right]$, where $\log(\cdot)$ is the binary
 380 logarithm and r_{al} and r_{bl} denote, as before, l -th element of \mathbf{r}_a and \mathbf{r}_b , respectively.
 381 It is important to note that KL is not symmetric. Hence, distance calculations using
 382 Δ_{KL} may result in non-symmetric distances.

383 *Chi-squared distance between profiles* (Chi2). A distance for categorical data,
 384 that has a strong link to the data visualization technique correspondence analysis,
 385 is the Chi-squared distance. There exist several forms and implementations of the
 386 chi-square distance that differ with respect to the chosen standardization. For an
 387 $n \times p$ contingency matrix $\mathbf{F}_j = \mathbf{Z}'_i \mathbf{Z}_j$, the Chi2 between rows a and b can be defined
 388 as $s \sum_{l=1}^p \frac{1}{f_{a\bullet}} \left(\frac{f_{al}}{f_{a\bullet}} - \frac{f_{bl}}{f_{b\bullet}} \right)^2$, where s is the sum of all elements of \mathbf{F} and \bullet denotes
 389 the summation in the appropriate dimension (rows or columns) of the matrix [see,
 390 e.g., 20, p.266]. In our notation, we can implement the Chi-squared distances
 391 as category dissimilarities by defining $\Phi(\mathbf{r}_a, \mathbf{r}_b) = \sum_{l=1}^{q_j} \frac{1}{p_l} (r_{al} - r_{bl})^2$, where we
 392 dropped the constant s , p_l corresponds to the l -th element of the j -th block of \mathbf{P}_d
 393 and, as before, r_{al} and r_{bl} denote the l -th element of \mathbf{r}_a and \mathbf{r}_b , respectively.

394 6. Supervised association-based distances

395 In a supervised setting, where we want to assign observations to classes (i.e.,
 396 categories) for one variable, say y , based on observations on categorical variables
 397 x_j where $j = 1, \dots, Q$, we can define a supervised variant of association-based
 398 categorical variable distances. That is, we can define category dissimilarities that
 399 take into account the association between variables y and x . Next, we can make
 400 predictions using either the K -nearest neighbors or a distance-based clustering
 401 method, where we fix the number of clusters to the number of classes and do a
 402 post-hoc comparison of clusters and classes. That is, we match the clusters to the
 403 true classes and assign labels accordingly. To define supervised association-based
 404 distances we create an $n \times c$ indicator matrix \mathbf{Z}_y , where c corresponds to the number
 405 of classes of y . If we add, to the right, this indicator matrix to \mathbf{Z} and insert this
 406 supplemented \mathbf{Z} into (1) through (2), we can calculate category dissimilarities using
 407 (4) and (5). Note that in this new setting, we have $Q + 1$ variables and consequently
 408 $Q + 1$ association-based category dissimilarity matrices $\mathbf{\Delta}^i$. However, the $(Q + 1)$ -
 409 th diagonal block gives the category dissimilarities between the categories of the y
 410 variable. In a supervised setting, a category (class) of y is to be predicted based on
 411 data from the other Q variables. The category dissimilarities for y should therefore
 412 not be used. This is easily achieved by simply ignoring these in the overall block
 413 diagonal category dissimilarity matrix $\mathbf{\Delta}$. That is, we construct $\mathbf{\Delta}$ by collecting
 414 only the first Q category dissimilarity matrices on its diagonal. As before, our
 415 framework allows for great flexibility in how to incorporate the information of
 416 variable y . In particular, the dissimilarity functions Φ^{ij} and the weights w_{ij} are

417 pair-specific. One could, as suggested in Section 5, set all weights w_{ij} equal to 1,
418 so that the category dissimilarities take into account all associations. We refer to
419 this choice as “full supervised” dissimilarity.

420 Alternatively, in a supervised setting, one may choose to have the category
421 dissimilarities depend only on the association with the variable y . This corresponds
422 to the choice $w_{ij} = 1$ for $j = Q + 1$ and 0 for all other pairs. In this case, category
423 dissimilarities may better discriminate with respect to the classes of y . We refer to
424 this choice as “supervised” dissimilarity. Note that both supervised variants lead to
425 new categorical variable distances. Moreover, as each of these require a choice of
426 association-based dissimilarity functions Φ^{ij} , for all pairs of variables, one could
427 consider these as a family of new, supervised categorical variable distances.

428 **7. Aggregation and dissimilarity tuning**

429 Our general framework introduced in Section 3 allows for great flexibility in the
430 implementation of distances between categorical variables. Indeed, as is clear from
431 Definition (5), all separate category dissimilarity measures can be combined and
432 aggregated according to the researcher’s preferences. How to exactly determine
433 which category dissimilarity and aggregation strategy is the most appropriate is
434 non-trivial, and this choice may depend on the properties of the data and the
435 research objectives.

436 In exploratory distance-based methods, e.g., cluster analysis, a clear measure
437 of fit is not available, since the goal is to find and interpret patterns in the data.
438 However, if a measure of fit can be calculated, it can be used to select an aggregation

439 and category-dissimilarity definition. That is, we can apply several aggregation
440 and category dissimilarity definitions, and compare the fit for each of them by
441 considering the selected measure. In a supervised classification setting, a choice
442 for aggregation and category dissimilarity definitions can then be made based
443 on the discrepancy between true classes with “predicted” classes. Therefore, the
444 aggregation state (i.e., w_{ij}) and the category dissimilarity function Φ^{ij} can be
445 treated as tuning parameters. Note, however, that 5 allows many combinations and
446 some choices need to be made to restrict the total search space.

447 8. Applications

448 To illustrate how our general framework can be used in practice, we consider
449 the following category dissimilarity measures: SM, Eskin, Lin, OF, Goodall 3,
450 Goodall 4, IOF, VE and VM (independent, 4), and, TVD, KL, Chi2, Supervised
451 TVD and Supervised TVD full (association-based, 5). In a supervised setting, we
452 employ K -nearest neighbors in which each new observation is labeled according
453 to a set of K *close* training points (neighbors). In an unsupervised setting, we
454 consider PAM as clustering method. For cluster analysis to find clusters that match
455 classes, one needs to assume that homogeneous groups in the data are related to
456 the class-variable. The supervised association-based distances do require classes
457 to be known as these are used for defining the distances. Consequently, a cluster
458 analysis using a supervised distance can no longer be considered to be unsupervised.
459 Furthermore, when considering matching of clusters to classes, it is expected that
460 supervised distances outperform the non-supervised distances.

461 We consider both synthetic and real-world datasets. The synthetic datasets
462 consist of 1000 observations described by 12 categorical variables with an underly-
463 ing cluster structure corresponding to 4 equal sized clusters. The datasets can be
464 characterized by mild and strong *clusterability* through the association of 4 of these
465 12 variables with a cluster membership variable. Hence, for these four variables,
466 the mild and strong *clusterability* levels correspond to variable-to-cluster Cramér's
467 V values of 0.5 and 0.7, respectively. The other eight variables serve as noise
468 variables. They are not associated with the underlying clusters. However, they are
469 related to each other. That is, they have pairwise associations among themselves
470 with a Cramér's V value of approximately 0.3. The two synthetic data sets are
471 available in the `cat_dist` package², and were generated using an evolutionary
472 algorithm. The complete data generating process is described in a step-by-step
473 fashion in [21].

474 Regarding the nine labeled real-world data sets in Table 1, all are available via
475 the UCI Machine Learning repository³. These data sets vary with respect to the
476 number of observations, variables and categories, and have been used in previous
477 papers [see, e.g., 13, 9] on categorical distances. For our analysis, each data set
478 is split into five folds for cross-validation. For the training subsets, we compute a
479 block diagonal matrix Δ , representing pairwise category dissimilarities for every
480 dissimilarity measure. The testing subset is used for performance assessment of

²The synthetic dataset names are `simcatdata1` for the mild structured dataset and `simcatdata2` for the strong structured dataset.

³<https://archive.ics.uci.edu/ml/index.php>

481 the considered methods. The chosen performance metrics are *i*) *Accuracy* for the
 482 nearest neighbors classifier, which represents the proportion of test observations
 483 that were correctly classified, and *ii*) the *Adjusted Rand Index* [ARI, 22], which
 484 measures the similarity between the assigned cluster of the test observations and
 485 their *true* cluster (i.e., the labels). These are common and well-documented choices
 486 for the specific supervised or unsupervised context, respectively. The procedure is
 487 iterated until each fold is used as test. The cross-validation procedure is repeated
 488 10 times.

Table 1: Data sets information

Data set	n	p	$\min q_j$	$\max q_j$	# clusters	Cramér's V
australian	690	8	2	14	2	0.26
balance	625	4	5	5	3	0.28
cars	1728	6	3	4	4	0.20
lymphography	148	18	2	8	4	0.38
soybean (large)	307	35	2	8	19	0.65
tae	151	5	2	46	3	0.30
tictactoe	958	9	3	3	2	0.14
vote	435	16	3	3	2	0.50
wbcd	699	9	9	11	2	0.77

489 8.1. *K-Nearest Neighbors classification of categorical data*

490 The KNN classification of the test observations is based on the calculation
 491 of the distance between each test observation and the training observations, as
 492 described in Section 3.1. Let $\mathbf{\Delta}_{train}$ denote the category dissimilarity matrix, where
 493 the subscript *train* indicates that if the category dissimilarities are data dependent,
 494 only observations of the training set were used. Furthermore, \mathbf{Z}_{test} and \mathbf{Z}_{train} are the

495 indicator matrices of the test and training observations, respectively. The distances
 496 of interest are in the columns of $\mathbf{Z}_{train}\Delta_{train}\mathbf{Z}'_{test}$ and the nearest neighbors for the
 497 j -th test observation are the K smallest values in the j -th column. The number
 498 of neighbors, K , is a hyper-parameter that can be tuned via cross-validation as
 499 described in Section 8.

500 To isolate the effect of the distance measure in the classification and clus-
 501 tering performance, we consider different values for K for each combination of
 502 distance metric and dataset. However, the choice of an optimal K is not cru-
 503 cial for our experiments and we limit ourselves to the following range of values
 504 $K \in \{1, 3, 7, 13, 21, 31, 43\}$.

505 Figure 1 presents the accuracy assessment of the tuned KNN classifier across
 506 the two scenarios and the category dissimilarity definitions. In particular, the left
 507 panel of the figure refers to mild variable-to-cluster associations, the right panel to
 508 strong variable-to-cluster associations. The size of each point is proportional to
 509 the chosen tuning hyper-parameter (i.e., the number of neighbors). The position of
 510 each point represents the median cross-validation accuracy across the 20 repetitions,
 511 and the error bars signify the interquartile range of the accuracy values across these
 512 repetitions. Within each data set or scenario, performance for the different category
 513 dissimilarity measures are ranked. In this way, the most effective measures in each
 514 scenario are immediately apparent as well as performances across scenarios.

515 We see that a strong structure yields higher accuracy, irrespective of the selected
 516 dissimilarity measure. Moreover, the Supervised TVD category dissimilarity (de-
 517 scribed in Section 6) consistently ranks as best performing measure. This is in line

518 with expectations as the structure in this study concerns the association strength of
 519 four of the variables with the clusters. The Supervised TVD category dissimilarity
 520 directly, and exclusively, incorporates these associations. The Supervised TVD
 521 full dissimilarity does not perform as well as Supervised TVD. This is because this
 522 measure also incorporates association among all the explanatory variables. As our
 523 data generating process also includes noise variables that are independent of the
 524 clusters, the pair-wise associations between them obscure the underlying cluster
 525 structure. This is also known as the *cluster masking problem* [see, for instance,
 526 21].

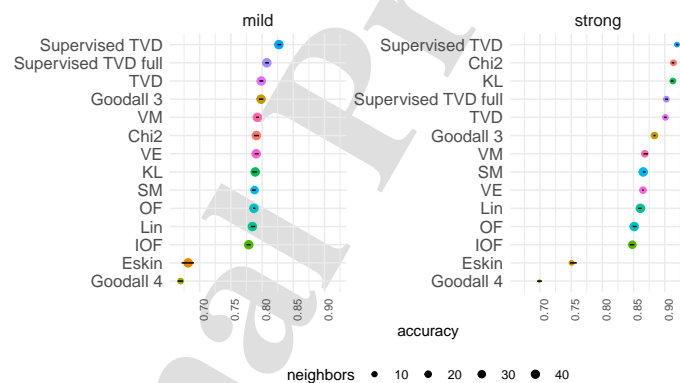


Figure 1: KNN classification accuracy for each considered distance measure and scenario

527 For the application on real data, we use the same range of values for K as
 528 for the experiments on synthetic data; however, we exclude the two largest levels
 529 (31 and 43) from due to the smaller size of some real-world datasets. For each
 530 data set/category dissimilarity combination, the value of K that minimizes the
 531 cross-validation estimate of the classifier's test accuracy is selected.

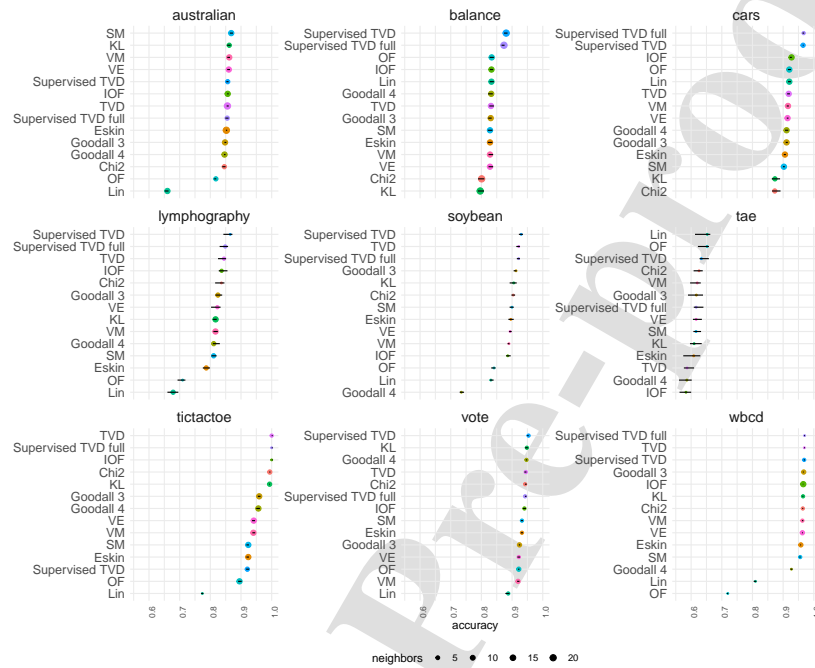


Figure 2: KNN classification accuracy for each considered distance measure and data set

532 Figure 2 depicts the accuracy assessment of the tuned KNN classifiers for each
 533 data set and category dissimilarity definition under study. The presentation format
 534 and the variability assessment of the results mirror those of Figure 1, with each
 535 panel in the figure corresponding to a specific real data set. For some data sets,
 536 e.g., *australian*, *vote*, and *wbcd*, the accuracy is high almost regardless of the
 537 chosen distance, with no variability over the 20 cross-validation repetitions. This
 538 indicates that in these three data sets, the considered categorical predictors have
 539 high discriminatory power for the underlying classes. Furthermore, since both the
 540 independent and the association-based measures perform well, taking into account

541 interactions (associations) between predictors is neither beneficial nor detrimental.

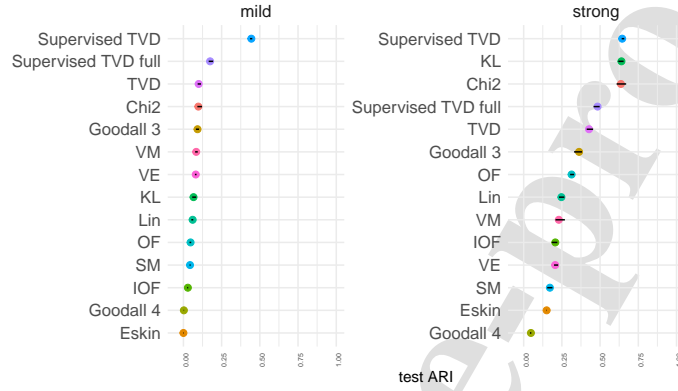


Figure 3: ARI results for PAM for each considered dissimilarity measure and scenario

542 8.2. Partitioning around medoids of categorical data

543 PAM is a distance-based iterative clustering procedure. Within a cluster, a
 544 medoid corresponds to the *median* observation, just like a centroid in K -means
 545 corresponds to the mean. The starting set of the K medoids is random, and each
 546 observation is assigned to the closest medoid; given the allocation of the obtained
 547 clusters, the medoids are updated accordingly. The procedure stops when there are
 548 no further changes in the set of medoids. Although we are working in an unsu-
 549 pervised context, the performance of PAM on each data set/category dissimilarity
 550 combination can be assessed via cross-validation by calculating distances based
 551 on the supervised association dissimilarities; this also allows for consistency with
 552 the KNN-based application. In particular, for each data set and each category
 553 dissimilarity definition, a medoids set is obtained by applying PAM to the train-
 554 ing data. That is, $\mathbf{D} = \mathbf{Z}_{train} \mathbf{\Delta}_{train} \mathbf{Z}'_{train}$ where for $\mathbf{\Delta}_{train}$ we consider all category

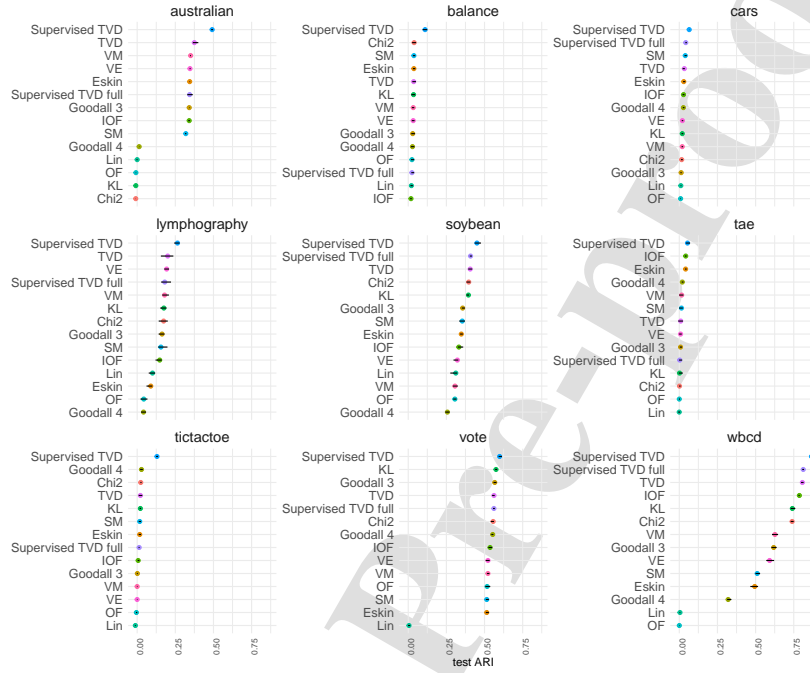


Figure 4: ARI results for PAM for each considered dissimilarity measure and data set

555 dissimilarity matrices described in Sections 4 and 5. Next, we compute the test
 556 observation-to-medoid distance matrix as $\mathbf{Z}_{medoid} \mathbf{\Delta}_{train} \mathbf{Z}'_{test}$, and assign each test
 557 observation to the cluster corresponding to the nearest medoid.

558 The PAM results are reported in Figures 3 (for the synthetic scenarios) and
 559 4 (for real data). These figures have similar characteristics to Figures 1 and 2,
 560 with a couple of distinctions: they present the ARI, which compares the clustering
 561 solution to the *true* cluster membership, and the size of the points is constant. As
 562 observed in the KNN-based application, a strong structure yields high ARI values,
 563 especially for the highest-ranking distance variants, which are all association-

564 based distances. Note that, when there is strong association, the supervised total
565 variance variant only performs slightly better than the unsupervised variants KL
566 and Chi2. The strong association of the four variables to the cluster structure
567 also leads to strong association among these variable that helps in picking up the
568 cluster structure. Furthermore, as the 8 noise variables are associated amongst
569 each other, this association obscures the association with the cluster structure so
570 that the supervised TVD variant only gives slightly better result than TVD, and
571 considerably worse results than the supervised variant that only considers the
572 association with the class variable. When there is only mild structure, cluster
573 retrieval is much worse. All unsupervised variants fail to adequately retrieve the
574 clusters. Regarding the supervised variants, we observe that the supervised TVD
575 variant that only uses the association of the variables with the class variable clearly
576 stands out. As before, the association among the noise variables, obscures the
577 association with the class variable.

578 The results on real data scenarios are reported in Figure 4. We observe that
579 the supervised association-based distance generally performs well. This is not a
580 surprise as, unlike for the unsupervised variants, the class labels are explicitly used
581 when determining the category dissimilarities and hence the distances. For the
582 *tae* data set, the PAM performance variability is very low, due to PAM failing to
583 find clusters that are related to the class variable. The cross-validated ARI values
584 obtained for each of the 20 random splits are zero, or close to zero, hence variation
585 is low. For slightly better results, e.g. using the Supervised TVD, higher ARI
586 values have an increased variability. In the *lymphography* data set case, the PAM

587 procedure does find some cluster structure (most ARI's are above .15), therefore
588 the 20 random splits of a small number of observations into five folds do affect the
589 cross-validated ARI values and cause the observed variability.

590 8.3. Summary of results

591 Our application of KNN and PAM to various data sets shows that supervised
592 and unsupervised performance is usually positively correlated. However, this
593 was not the case for some datasets like *balance*, *lymphography*, and *tictactoe*,
594 where clusters did not match classes well. Using supervised distance slightly
595 improved results. On synthetic data, association-based measures are preferred for
596 clustering related to classes. Real datasets with high ARI values also benefited
597 from association-based distances, such as in the *wbcd* dataset where four out of
598 five measures with ARI value exceeding 0.7 were association-based.

599 9. Conclusion

600 In this paper, we introduced a comprehensive framework for implementing
601 distances between categorical variables that is flexible, efficient, and transparent.
602 We showed that both independent and association-based distances can be integrated
603 within our framework. Consequently, our approach can be used to apply existing
604 measures, thus enhancing implementation transparency, or to introduce new, highly
605 customizable ones.

606 Our framework is not limited to specific methods or applications, allowing
607 for domain-specific dissimilarity measures. Specifically in supervised contexts,

608 our approach accommodates measures considering associations with response
609 variables. The dissimilarity definitions for independent categories outlined in
610 Section 4 are integrated into the `nomclust` R package [12], excluding ordered
611 category dissimilarities. The package does not support association-based measures.
612 Conversely, the `catdist` package ⁴ implements both the independent and the
613 association-based measures discussed in this paper.

614 To underscore the significance of selecting the “optimal” distance for a given
615 problem, we employed our framework in both supervised (KNN) and unsupervised
616 (PAM) settings. For real-world datasets, the choice of measure did not significantly
617 affect results unless variables lacked discriminatory power or strong clustering,
618 or when KNN/PAM were not suitable methods for the problem. In some scenar-
619 ios, all measures delivered comparable results. However, setting aside extreme
620 cases, opting for the “most appropriate” measure often improved outcomes, with
621 association-based measures generally performing better than independent category
622 measures. While further validation with more datasets would be beneficial, it falls
623 outside the scope of this paper.

624 Our framework allows for implementing new or customized distances easily,
625 as demonstrated by the supervised total variation distances introduced in Section
626 6. Instead of introducing and evaluating new measures, a more valuable approach
627 would be to systematically compare existing dissimilarity measures for categorical
628 variables in the literature, highlighting their strengths and weaknesses.

⁴Available on GitHub at [https://github.com/alfonsoIodiceDE/catdist_](https://github.com/alfonsoIodiceDE/catdist_package)
package.

629 Finally, while our focus is on categorical variables, our framework can be
630 adapted for mixed-variable settings with numerical variable distances. This inte-
631 gration is not a trivial task and requires addressing scale differences within and
632 between variable types, yet our framework offers a solid starting point for this task.

633 **Funding**

634 This research did not receive any specific grant from funding agencies in the
635 public, commercial, or not-for-profit sectors.

636 **Conflict of interest**

637 The authors declare that they have no conflict of interest.

638 **Data availability**

639 Extended results and the code to reproduce them are available online at `https :`
640 `//alfonsoiodicede.github.io/blogposts_archive/distances_`
641 `experiment_superv_unsuperv.html`. The data used in this study are
642 available in the `catdist` package⁵.

643 **References**

644 [1] E. Blanco-Mallo, L. Morán-Fernández, B. Remeseiro, V. Bolón-Canedo, Do
645 all roads lead to Rome? Studying distance measures in the context of machine
646 learning, *Pattern Recognition* 141 (2023) 109646.

⁵https://github.com/alfonsoIodiceDE/catdist_package

- 647 [2] G. James, D. Witten, T. Hastie, R. Tibshirani, J. Taylor, *An Introduction to*
648 *Statistical Learning: With Applications in Python*, Springer Nature, 2023.
- 649 [3] L. Kaufman, P. Rousseeuw, *Finding groups in data: an introduction to cluster*
650 *analysis*, John Wiley & Sons, New York, 1990.
- 651 [4] I. Borg, P. Groenen, *Modern multidimensional scaling: Theory and applica-*
652 *tions*, Springer Science & Business Media, 2005.
- 653 [5] J. C. Gower, S. G. Lubbe, N. J. Le Roux, *Understanding Biplots*, John Wiley
654 & Sons, 2011.
- 655 [6] S. Le, T. Ho, An association-based dissimilarity measure for categorical data,
656 *Pattern Recognition Letters* 26 (2005) 2549–2557.
- 657 [7] A. Ahmad, L. Dey, A k-mean clustering algorithm for mixed numeric and
658 categorical data, *Data & Knowledge Engineering* 63 (2007) 503–527.
- 659 [8] H. Jia, Y. Cheung, J. Liu, A new distance metric for unsupervised learning
660 of categorical data, *IEEE Transactions on Neural Networks and Learning*
661 *Systems* 27 (2014) 1065–1079.
- 662 [9] M. Ring, F. Otto, M. Becker, T. Niebler, D. Landes, A. Hotho, *Condist: A*
663 *context-driven categorical distance measure*, in: A. Appice, P. Rodrigues,
664 V. Santos Costa, C. Soares, J. Gama, A. Jorge (Eds.), *Machine Learning and*
665 *Knowledge Discovery in Databases*, Springer International Publishing, Cham,
666 2015, pp. 251–266.

- 667 [10] E. Mousavi, M. Sehati, A generalized multi-aspect distance metric for
668 mixed-type data clustering, *Pattern Recognition* 138 (2023) 109353.
- 669 [11] H. Rezaei, N. Daneshpour, Mixed data clustering based on a number of
670 similar features, *Pattern Recognition* 143 (2023) 109815.
- 671 [12] Z. Šulc, H. Řezanková, Comparison of similarity measures for categorical
672 data in hierarchical clustering, *Journal of Classification* 36 (2019) 58–72.
- 673 [13] S. Boriah, V. Chandola, V. Kumar, Similarity measures for categorical data:
674 A comparative evaluation, in: *Proceedings of the 2008 SIAM International
675 Conference on Data Mining*, SIAM, 2008, pp. 243–254.
- 676 [14] E. Eskin, A. Arnold, M. Prerau, L. Portnoy, S. Stolfo, A geometric framework
677 for unsupervised anomaly detection, in: *Applications of data mining in
678 computer security*, Springer, 2002, pp. 77–101.
- 679 [15] D. Lin, An information-theoretic definition of similarity, in: *Proceedings
680 of the Fifteenth International Conference on Machine Learning*, Morgan
681 Kaufmann Publishers Inc., San Francisco, CA, USA, 1998, p. 296–304.
- 682 [16] K. Spärck Jones, A statistical interpretation of term specificity and its appli-
683 cation in retrieval, *Journal of Documentation* 28 (1972) 11–21.
- 684 [17] D. Goodall, A new similarity index based on probability, *Biometrics* (1966)
685 882–907.

- 686 [18] H. Drost, Philentropy: information theory and distance quantification with R,
687 Journal of Open Source Software 3 (2018) 765.
- 688 [19] S. Kullback, R. Leibler, On information and sufficiency, The Annals of
689 Mathematical Statistics 22 (1951) 79–86.
- 690 [20] A. Gifi, Nonlinear multivariate analysis, John Wiley & Sons Ltd., 1990.
- 691 [21] M. van de Velden, A. Iodice D’Enza, F. Palumbo, Cluster correspondence
692 analysis, Psychometrika 82 (2017) 158–185.
- 693 [22] L. Hubert, P. Arabie, Comparing partitions, Journal of Classification 2 (1985)
694 193–218.

A general framework for implementing distances for categorical variables

Michel van de Velden^a, Alfonso Iodice D'Enza^b, Angelos Markos^c, Carlo Cavicchia^d

^a*Econometric Institute, Erasmus University Rotterdam*

^b*Department of Political Sciences, University of Naples Federico II*

^c*Department of Primary Education, Democritus University of Thrace*

^d*Econometric Institute, Erasmus University Rotterdam*

-
-
- The quantification of distances among objects plays an important role in many statistical methods.
 - An appropriate definition of a distance depends on the nature of the data and the problem at hand.
 - For categorical data, the definition of a distance is complex, as there is no straightforward quantification of the size of the observed differences.
 - A general framework that allows for an efficient and transparent implementation of distances between observations on categorical variables is needed.
 - In a supervised classification setting, the framework can be used to construct distances that incorporate the association between the response and predictor variables and hence improve the performance of distance-based classifiers.

A general framework for implementing distances for categorical variables

Michel van de Velden^a, Alfonso Iodice D'Enza^b, Angelos Markos^c, Carlo Cavicchia^d

^a*Econometric Institute, Erasmus University Rotterdam*

^b*Department of Political Sciences, University of Naples Federico II*

^c*Department of Primary Education, Democritus University of Thrace*

^d*Econometric Institute, Erasmus University Rotterdam*

-
-
- Michel van de Velden is an associate Professor of Statistics at the Econometric Institute of the Erasmus University Rotterdam. His main research interests are exploratory data analysis. In particular, dimension reduction and cluster analysis methods with a strong focus on data visualization. In addition, he is involved in several supervised machine learning projects involving tree-based machine learning methods.
 - Alfonso Iodice D'Enza is an associate Professor of Statistics at the University of Naples Federico II (Italy). His areas of interest include statistical learning, clustering, dimension reduction, computational statistics and visualisation, with applications in behavioural sciences.
 - Angelos Markos is an associate Professor of Data Analysis in the Social Sciences at the Democritus University of Thrace, School of Education (Greece). His areas of interest include multivariate analysis (especially dimension reduction and clustering), psychological testing and measurement in the social sciences and statistics education.
 - Carlo Cavicchia is an assistant Professor of Statistics at the Econometric Institute of the Erasmus University Rotterdam. His research is focused on the methodological and computational aspects of data analysis. He is currently working on latent variable models, unsupervised classification and model-based composite indicators.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre