


DataSimilarity: Quantifying Similarity of Datasets and Multivariate Two- and k -Sample Testing

Marieke Stolte 
TU Dortmund University

Luca Sauer
TU Dortmund University

Jörg Rahnenführer 
TU Dortmund University

Andrea Bommert 
TU Dortmund University

Abstract

Quantifying the similarity of two or more datasets is a common task in many applications in statistics and machine learning, such as two- or k -sample testing and meta- or transfer learning. The **DataSimilarity** package contains a variety of methods for quantifying the similarity of datasets. The package includes 36 methods of which 14 are implemented for the first time in R. The remaining are wrapper functions for methods with already existing implementations that unify and simplify the various input and output formats of different methods and bundle the methods of many existing R packages in a single package. Here, we give details on the methods and implementations and show some application examples.

Keywords: dataset similarity, two-sample testing, multi-sample testing.

1. Methods

In the following, we describe the general setup in the two- or k -sample problem that most of the implemented methods have in common. Moreover, we discuss the selection of the implemented methods and present one example method for each application domain in more detail.

1.1. The two- and k -sample problem

Most methods for quantifying the similarity of datasets are proposed in the literature as test statistics for two- or k -sample testing. For this, a dataset is seen as a sample from a set of random variables that follow some true underlying distribution. Often, the similarity or distance of these underlying distributions is estimated.

In the following, we assume that at least two different datasets $X^{(1)}$ and $X^{(2)}$ are given consisting of n_1 and n_2 samples $X_1^{(1)}, \dots, X_{n_1}^{(1)} \sim F_1$ and $X_1^{(2)}, \dots, X_{n_2}^{(2)} \sim F_2$, respectively. We assume $X_i^{(1)}, X_j^{(2)} : \mathcal{X} \rightarrow \mathbb{R}^p \forall i \in \{1, \dots, n_1\}, j \in \{1, \dots, n_2\}$ and call the p components of each sample features or variables. The two-sample problem is defined as the testing problem

$$H_0 : F_1 = F_2 \text{ vs. } H_1 : F_1 \neq F_2. \quad (1)$$

This testing problem is sometimes also called testing for homogeneity of the two distributions. In some cases, it is assumed that there are n_i observations of a target variable Y in each dataset. However, most methods only require the feature variables and cannot deal with a target variable in a meaningful way.

Analogously to the two-sample problem, the k -sample or multi-sample problem is defined for $k \geq 2, k \in \mathbb{N}$, datasets $X^{(1)}, \dots, X^{(k)}$ with sample sizes $n_i, i = 1, \dots, k$, as

$$H_0 : F_1 = F_2 = \dots = F_k \text{ vs. } H_1 : \exists i \neq j \in \{1, \dots, k\} : F_i \neq F_j,$$

where F_i denotes the distribution of each sample in the i th dataset.

Each of the considered methods can be seen as a measure of similarity or distance between the $F_i, i = 1, \dots, k$. Not all of these methods include a hypothesis test.

We use the hat symbol to denote estimators. We denote the pooled sample as $\{Z_1, \dots, Z_N\} = \{X_1^{(1)}, \dots, X_{n_1}^{(1)}, \dots, X_1^{(k)}, \dots, X_{n_k}^{(k)}\}$, where $N = \sum_{i=1}^k n_i$ is the total sample size. Additionally, we assume that all Z_i are distributed independently.

1.2. Selection of methods

Previously, in a comprehensive literature review ([Stolte, Kappenberg, Rahnenführer, and Bommert 2024](#)), 118 methods were described and divided into the ten classes

1. Comparison of cumulative distribution functions, density functions, or characteristic functions,
2. Methods based on multivariate ranks,
3. Discrepancy measures for distributions,
4. Graph-based methods,
5. Methods based on inter-point distances,
6. Kernel-based methods,
7. Methods based on binary classification,
8. Distance and similarity measures for datasets,
9. Comparison based on summary statistics, and
10. Testing approaches.

Moreover, the methods were compared with respect to 22 criteria judging their applicability, interpretability, and theoretical properties. The **DataSimilarity** package comprises 36 methods that fulfill at least one of the following properties:

1. The method is implemented in R.
2. The method is one of the top methods ordered by the highest number of fulfilled criteria, and fulfills at least 11 criteria of the 20 criteria, excluding the consistency criteria.

3. The method is the best in its subclass in the theoretical comparison, and no other method from this subclass was chosen based on the first two criteria.

To avoid preferring methods that define a test over methods that do not, and therefore can by definition not fulfill the consistency criteria, consistency is not counted for determining the top methods. We chose 11 as the cutoff for the number of fulfilled criteria, as this is the range where the implemented methods typically lie, and it ensures that at least more than half of the criteria are fulfilled.

1.3. Definition of example methods

In the following, we differentiate six cases with regard to the applicability of the selected methods. These are summarized in Table 1. We always indicate which method is applicable in which case. In the following, we explain one example method for each case. These methods are used later in examples for applying the **DataSimilarity** package. Brief descriptions of the remaining methods can be found in Section 5.

Scenario no.	No.datasets	Scale level	Target variable
1	$k = 2$	Numeric	No
2	$k \geq 2$	Numeric	No
3	$k = 2$	Numeric	Yes
4	$k = 2$	Categorical	No
5	$k \geq 2$	Categorical	No
6	$k = 2$	Categorical	Yes

Table 1: Overview of considered cases for applicability of the dataset similarity methods. If present, the target variable included in each dataset has to be a categorical variable.

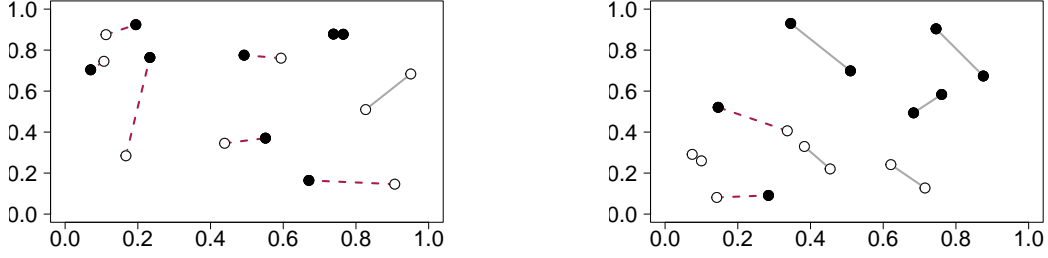
1. Methods applicable to exactly two numeric datasets without target variables

One example method for this case is the [Rosenbaum \(2005\)](#) cross-match test. It is a graph-based method. Most graph-based methods work by constructing a similarity graph on the pooled sample and counting the edges that connect points from different samples. Here, the optimal non-bipartite matching is used, i.e., a graph where pairs of two observations in the pooled sample are connected such that the sum over the edge lengths (= Euclidean distances of connected observations) is minimized. The optimal non-bipartite matching for two example data situations is shown in Figure 1. In case of an odd number of observations, a ghost observation is introduced that has the highest distance to all other observations. The observation that is matched with that ghost observation is discarded from further analysis.

The test statistic of the cross-match test is given by the standardized cross-match count

$$\frac{\text{CMC} - \mathbb{E}_{H_0}(\text{CMC})}{\sqrt{\text{VAR}_{H_0}(\text{CMC})}},$$

where CMC denotes the cross-match count and \mathbb{E}_{H_0} and VAR_{H_0} its expectation and variance, respectively, under $H_0 : F_1 = F_2$. The cross-match count is the number of edges connecting



(a) Datasets drawn from the same distribution. (b) Datasets drawn from different distributions.

Figure 1: Optimal non-bipartite matching for example datasets. Dataset 1 is indicated by white points and Dataset 2 by black points. Edges connecting points from different datasets are indicated by red and dashed lines. Edges connecting points from the same sample are indicated by grey and solid lines.

points stemming from different datasets. The exact distribution of the test statistic under H_0 is known. For small samples, it can be used for computing an exact p value. For large samples, the asymptotic standard normal distribution of the test statistic can be used. The idea of the test is that for similar datasets, the number of edges connecting points from different samples is expected to be higher than in datasets that differ. This is illustrated in Figure 1a compared to Figure 1b. In case of data drawn from different datasets, fewer edges connect points from different datasets, indicated by the lower number of red edges in Figure 1b.

2. Methods applicable to two or more numeric datasets without target variables.

The method of Mukherjee, Agarwal, Zhang, and Bhattacharya (2022) is an extension of the Rosenbaum (2005) cross-match test for multiple samples. The cross-match counts $A = (a_{12}, a_{13}, \dots, a_{ik}, a_{23}, \dots, a_{2k}, \dots, a_{k-1,k})^\top$ for all pairs of datasets are calculated using the optimal non-bipartite matching on the pooled sample. The test statistic then is the Mahalanobis distance of the observed cross-counts under the null hypothesis $H_0 : F_1 = F_2 = \dots = F_k$

$$\text{MMCM} = (A - E_{H_0}(A))^\top \text{COV}_{H_0}^{-1}(A)(A - E_{H_0}(A)).$$

The expectation and covariance matrix of the cross-count vector A under H_0 can be calculated analytically and depend only on the sample sizes $n_i, i = 1, \dots, k$. Small values of the multi-sample Mahalanobis cross-match (MMCM) statistic indicate similarity. However, as there is no known upperbound, it is hard to interpret the MMCM value. The MMCM statistic follows a $\chi^2_{\binom{k}{2}}$ distribution asymptotically under the null, which can be used for testing.

3. Methods applicable to exactly two numeric datasets with target variables

Ntoutsis, Kalousis, and Theodoridis (2008) propose measuring dataset similarity based on probability density estimates derived from decision trees. For this, it is assumed that in addition to both covariate datasets $X^{(1)}$ and $X^{(2)}$, categorical target variables $Y^{(1)}$ and $Y^{(2)}$ are given. On each dataset $X^{(j)}$, a classification tree is constructed with $Y^{(j)}$ as the target variable, $j = 1, 2$. The splits defined by the decision trees induce a partition of the feature space \mathcal{X} such that each leaf node corresponds to one segment in the partition. Figure 2 demonstrates the procedure for two example datasets. First, trees are fit to each dataset

(Figure 2a and 2b). Then, the sample space is divided into segments based on the splits performed in each tree (Figure 2c and 2d). These partitions are intersected (Figure 2e) and based on the joint partition, the probability densities $P_D(\mathcal{X})$ and $P_D(Y^{(j)}, \mathcal{X})$ are estimated for $D \in \{X^{(1)}, X^{(2)}, Z\}$.

Let n_r denote the number of segments in the joint partition and n_c the number of classes in $X^{(1)}$ and $X^{(2)}$. $\hat{P}_D(\mathcal{X}) \in \mathbb{R}^{n_r}$ uses the proportion of observations in D that fall into each segment of the joint partition. This means that for each of the n_r segments of the partition, the number of observations from dataset D that fall into that segment is counted and divided by the total number of observations in D . For the estimation of the joint density $P_D(Y, \mathcal{X})$, the proportion of observations that fall into each segment of the joint partition and belong to each class is determined, $\hat{P}_D(Y, \mathcal{X}) \in \mathbb{R}^{n_r \times n_c}$. Here, for each of the n_r segments of the partition and each of the n_c classes, the number of observations in D where the corresponding target variable has the respective class value and that fall into the respective segment is counted and divided by the total number of observations in D . The conditional density $P_D(Y|\mathcal{X})$ is estimated by calculating the proportion of observations belonging to each class separately for each segment, $\hat{P}_D(Y|\mathcal{X}) \in \mathbb{R}^{n_r \times n_c}$. Here, for each of the n_r segments of the partition and each of the n_c classes, the number of observations in D where the corresponding target variable has the respective class value and that fall into the respective segment is counted and divided by the total number of observations in D that fall into the respective segment.

Then, Ntoutsi *et al.* (2008) consider the similarity index

$$s(p, q) = \sum_i \sqrt{p_i \cdot q_i}$$

for vectors p and q , where $(n_r \times n_c)$ -matrices are interpreted as $(n_r \cdot n_c)$ -dimensional vectors. For the conditional distribution, the similarity vector $S(Y|\mathcal{X}) \in \mathbb{R}^{n_r}$ is computed with $S(Y|\mathcal{X})_i = s(\hat{P}_{X^{(1)}}(Y|\mathcal{X})_{i\bullet}, \hat{P}_{X^{(2)}}(Y|\mathcal{X})_{i\bullet})$ and index $i\bullet$ denoting the i -th row. Based on this, three similarity measures for datasets are proposed:

1. NTO1 = $s(\hat{P}_{X^{(1)}}(\mathcal{X}), \hat{P}_{X^{(2)}}(\mathcal{X}))$
2. NTO2 = $s(\hat{P}_{X^{(1)}}(Y, \mathcal{X}), \hat{P}_{X^{(2)}}(Y, \mathcal{X}))$
3. NTO3 = $S(Y|\mathcal{X})^\top \hat{P}_Z(\mathcal{X})$.

All three measures have values in the interval $[0, 1]$, where high values correspond to high similarity.

4. Methods applicable to exactly two categorical datasets without target variables

Hediger, Michel, and Näf (2022) provide a two-sample test based on random forests that is applicable for both numeric and categorical data. For this, a pooled dataset is created where each observation is labeled according to its original dataset membership, and a random forest is trained to distinguish between the dataset labels. The idea is that if the datasets are generated from the same distribution, the classification error of the random forest should be close to the chance level, otherwise, the classifier should be able to distinguish between the two distributions and hence the classification error should be lower than the chance level. One advantage of using random forests as the classifier is that it requires almost no tuning.

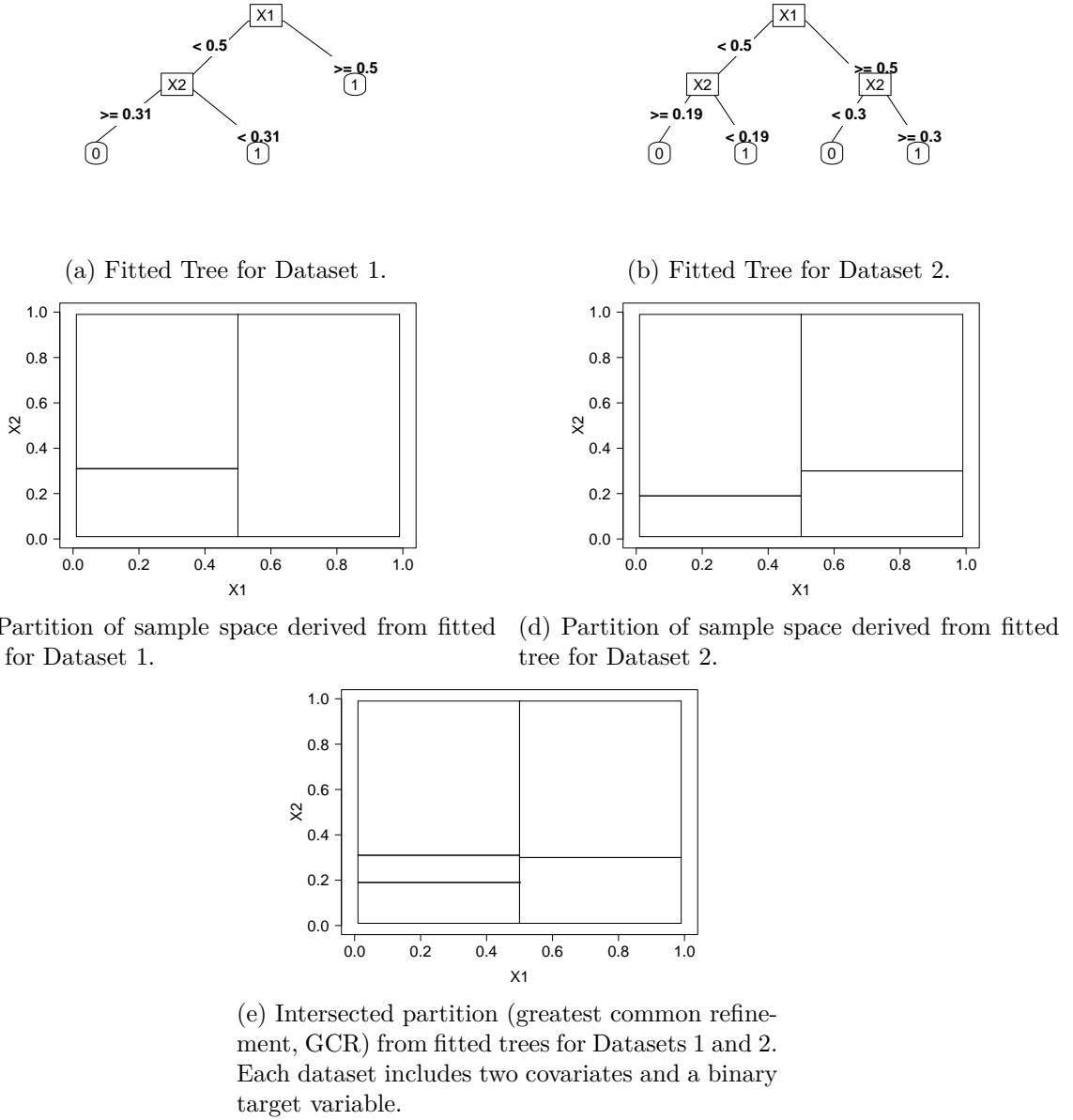


Figure 2: Partitioning of sample space by fitting trees to two example datasets.

An asymptotic test is proposed. For this, the pooled dataset has to be split into a training set on which the random forest is trained and a test set on which its classification error is evaluated. In the implementation, both datasets are split in half to create a training and a test dataset. Alternatively, an out-of-bag (OOB) based permutation test can be performed that does not require data splitting. OOB statistics can be used to increase the sample efficiency compared to the test based on a holdout sample. Both the OOB-based test and the asymptotic version of the test are implemented. The test statistic is either the mean of the per-class OOB or test classification errors or the overall OOB or test classification

error, respectively. In the asymptotic case, a binomial test is performed in case of the overall classification error, or a Z test is performed in case of the mean per-class classification error. Otherwise, a permutation test is performed. The variable importance measures of the random forest can provide additional insights into sources of distributional differences.

5. Methods applicable to two or more categorical datasets without target variables

The general idea of [Lopez-Paz and Oquab \(2017\)](#) is to use a classifier to determine which of two or more datasets a sample belongs to. The *classifier two-sample test (C2ST)* uses the classification accuracy of this classifier as its test statistic.

The C2ST consists of five steps:

1. Construct the dataset consisting of the samples from all datasets, labeled with their membership to each of the datasets.
2. Assign the observations of the dataset constructed in 1. randomly to a training and test set.
3. Train a classifier that predicts for an observation to which dataset $X^{(j)}, j = 1, \dots, k$ it belongs.
4. Calculate the C2ST statistic, which is the accuracy on the test set. The accuracy should be close to the chance level for $F_1 = \dots = F_k$, and it should be greater than the chance level for $\exists i \neq j \in \{1, \dots, k\} : F_i \neq F_j$ since in the latter case the classifier should identify distributional differences between the samples.
5. Calculate a p value using a binomial test for comparing the accuracy to the chance level.

Maximizing the power of a C2ST is a trade-off between using a large training set to optimize the classifier and a large test set to better evaluate the performance of the classifier.

The test statistic is interpretable as the percentage of samples that are correctly classified on the unseen test data. The above-mentioned test of [Hediger et al. \(2022\)](#) can be seen as a special case of the general framework proposed by [Lopez-Paz and Oquab \(2017\)](#). One difference in the implementation of the tests is that for the C2ST, categorical data is dummy-encoded, while for the test of [Hediger et al. \(2022\)](#) the categorical variables are passed to `ranger::ranger()` directly. Moreover, the use of OOB predictions and feature importance is specific to the random forest-based test and cannot be used for all of the available classifiers for the C2ST. Further, the C2ST uses the accuracy as its test statistic while the test of [Hediger et al. \(2022\)](#) uses the classification error, i.e., $1 - \text{accuracy}$.

6. Methods applicable to exactly two categorical datasets with target variables

[Alvarez-Melis and Fusi \(2020\)](#) define a distance based on optimal transport between datasets that include a target (class) variable Y . The *optimal transport dataset distance (OTDD)* is defined as

$$d_{\text{OT}}(X^{(1)}, X^{(2)}) = \min_{\pi \in \Pi(F_1, F_2)} \int_{\mathcal{Z} \times \mathcal{Z}} d_{\mathcal{Z}}(z, z')^q d\pi(z, z')$$

where $X^{(1)}, X^{(2)}$ denote the two datasets,

$$\Pi(F_1, F_2) := \{\pi_{1,2} \in \mathcal{P}(\mathcal{Z} \times \mathcal{Z}) | \pi_1 = F_1, \pi_2 = F_2\}$$

is the set of joint distributions over the product space $\mathcal{Z} \times \mathcal{Z}$ over the sample space of the pooled sample with marginal distributions F_1 and F_2 , and

$$d_{\mathcal{Z}}(z, z') := (d_{\mathcal{X}}(x, x')^q + W_{q'}^{q'}(\alpha_y, \alpha_{y'}))^{1/q}.$$

defines a distance of two points $z^\top = (x^\top, y)$, and $z'^\top = (x'^\top, y')$ in the pooled sample. $d_{\mathcal{X}}$ defines a distance on the covariate space, e.g., the Euclidean distance, and $W_{q'}(\alpha_y, \alpha_{y'})$ is the q' -Wasserstein distance of the distribution of the subset of covariate data with corresponding response value y and the distribution of the subset of covariate data with corresponding response value y' . The powers q and q' have to be chosen in advance to calculate the OTDD. The optimal transport problem can intuitively be motivated by imagining each probability density as a pile of dirt. Then, the cost function corresponds to the cost for transporting the dirt from one point to another, which is proportional to the distance between the two points. The optimal transport then corresponds to the lowest cost required for moving one pile of dirt fully to the shape and location of the other. Therefore, distributions can be regarded as more similar if the optimal transport between them is lower. For an intuitive explanation and visualization of the OTDD, also refer to [Alvarez-Melis and Fusi \(2020\)](#).

2. General comments on implementation

Where possible, existing implementations are used. If methods already have a name in the article where they were proposed or in the secondary literature, the corresponding functions are named after that, e.g., `Wasserstein()` for the Wasserstein distance, `MMD()` for the maximum mean discrepancy (MMD), or `CMDistance()` for the constrained minimum (CM) distance. Otherwise, the function names are composed of the first letters of the surnames of all authors of the article where the respective method was originally proposed, e.g., `FR()` for the Friedman-Rafsky test proposed by [Friedman and Rafsky \(1979\)](#), or the full surname in case of a single author, e.g., `Bahr()` for the test proposed by [Bahr \(1996\)](#). The input and output of the methods from different existing packages and of the newly implemented methods are unified. To achieve this, for some existing methods, it was sufficient to implement a wrapper calling the original function.

In other cases, we re-implemented the method from scratch if the R package was archived and additional issues with the original implementation occurred. This was the case for the DiPro-Perm test ([Wei, Lee, Wichers, and Marron 2016](#)) for which the original implementation in the **diproperm** package ([Allmon, Marron, and Hudgens 2021](#)) yields non-reproducible results. Moreover, the implementations of the multi-sample cross-match test of [Petrie \(2016\)](#) and the previously mentioned multi-sample Mahalanobis cross-match test (MMCM) of [Mukherjee et al. \(2022\)](#) in the **multicross** package ([Agarwal, Bhattacharya, and Zhang 2020](#)) could not be used due to the output format that made it impossible to access the test statistic and p value. More details on the new implementations compared to the aforementioned versions can be found in Section 5.

Each method gets two (or more) datasets as its first input parameters. After that, arguments specific to the method follow. E.g., many methods perform a permutation test for which the number of permutations (`n.perm`) has to be specified. The output is of class `'htest'` and includes

- `statistic`: The test statistic

- **parameter** (optional): A parameter specifying the null distribution (e.g., degrees of freedom for a χ^2 distribution).
- **p.value**: The p value (if an asymptotic or permutation / Bootstrap test is performed).
- **estimate**: The sample estimate(s) (if available, e.g., the edge count for edge-count tests, NULL for many methods).
- **alternative**: The alternative hypothesis. For two datasets, this is $F_1 \neq F_2$, for k datasets it is $\exists i \neq j \in \{1, \dots, k\} : F_i \neq F_j$.
- **data.name**: Names of the supplied datasets.
- Further elements specific to the method (optional), e.g., the variable importances for the test of Hediger *et al.* (2022).

We use the ‘**htest**’ class as it is widely adopted for storing results of hypothesis tests in R and most of the implemented methods are two- or k -sample tests. Objects of class ‘**htest**’ will be automatically printed in an appealing format using the `print.htest()` function from the **stats** package. For methods for which no test is performed, the **p.value** is set to NULL. This allows pretty printing of the results and a unified output format for the corresponding functions. For many of the newly implemented permutation tests, we use the `boot()` function from the **boot** package that is included in R.

In typical applications, users should choose a test a priori and not based on test results. Therefore, the new functions perform exactly one test and return only the results corresponding to that single test. Some of the former implementations used to perform multiple tests based on the same metrics or always returned the asymptotic p value in addition to a permutation p value. This could lead to unscientific practices like choosing the test based on the desired result. As an exception, for implementations that output multiple related tests, we offer wrapper functions that also perform these multiple tests. Often, conducting them at once is computationally faster than performing each test individually when large parts of the calculation are the same. This option might be useful in certain situations where multiple tests need to be applied to the same data, e.g., when performing method comparison studies. We do not advise applying multiple tests for the same hypothesis on the same datasets when conducting inference for a specific real-life application.

Some of the existing implementations already include setting a random seed, and some do not. Therefore, for unity, the new methods all include a random seed argument and set the random seed to the supplied value for reproducibility.

3. Illustrations

In the following, the example methods for the six cases from Section 1.3 are applied to some real-world datasets. These are typically subsets of a dataset defined in such a way that, from the application background, it is clear that the subsets should or should not differ. The datasets were selected from the datasets included in the R packages that the **DataSimilarity** package depends on, so no additional packages are needed. To apply all the methods, we simply need to load the **DataSimilarity** package.

```
R> library("DataSimilarity")
```

3.1. Exactly two numeric datasets without target variables

The dataset `dhfr` (Sutherland and Weaver 2004) from the `caret` package (Kuhn and Max 2008) is a binary classification dataset (regarding Dihydrofolate Reductase inhibition) consisting of 325 compounds of which 203 are labeled as ‘active’ and 122 as ‘inactive’. The variables are 228 molecular descriptors. As the active and inactive compounds should differ in their descriptors, we divide the dataset according to the first variable that indicates the activity status.

```
R> data(dhfr, package = "caret")
R> act <- dhfr[dhfr$Y == "active", -1]
R> inact <- dhfr[dhfr$Y == "inactive", -1]
```

We apply the Rosenbaum cross-match test to check whether the active and inactive compounds differ. As the combined sample size is smaller than 340, we can apply the exact test:

```
R> Rosenbaum(act, inact, exact = TRUE)
```

```
Exact cross-match test
```

```
data: act and inact
z = -9.4098, p-value < 2.2e-16
alternative hypothesis: The distributions of act and inact are unequal.
sample estimates:
edge.count
      20
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(act, inact, method = "Rosenbaum", exact = TRUE)
```

```
Exact cross-match test
```

```
data: act and inact
z = -9.4098, p-value < 2.2e-16
alternative hypothesis: The distributions of act and inact are unequal.
sample estimates:
edge.count
      20
```

The cross-match count is equal to 20. At most, there could be 122 cross-matches if each observation from the ‘inactive’ dataset was connected to an observation in the ‘active’ dataset. Therefore, the cross-match count of 20 can be considered a rather small value. This is also reflected by the z score of -9.41. Consequently, we see that the hypothesis of equal distributions can be rejected with a p value smaller than $2.2 \cdot 10^{-16}$.

We obtain a warning that informs us that a ghost value was introduced when calculating the optimal non-bipartite matching, due to the odd pooled sample size. This means that an artificial point was added to the sample that has the highest distance to all other points in the sample, such that the optimal non-bipartite matching, which needs an even sample size, could be calculated. The ghost value and the point with which it was matched are then discarded from the subsequent calculations.

3.2. More than two numeric datasets without target variables

The well-known `iris` dataset (Fisher 1936) included in the `datasets` package that comes with base R (R Core Team 2024) includes measurements of sepal and petals of 50 flowers each of three iris species. We compare the datasets for the three species *Iris setosa*, *versicolor*, and *virginica*.

```
R> data("iris")
R> setosa <- iris[iris$Species == "setosa", -5]
R> versicolor <- iris[iris$Species == "versicolor", -5]
R> virginica <- iris[iris$Species == "virginica", -5]
```

For comparing the three datasets, we use the Mukherjee *et al.* (2022) Mahalanobis multisample crossmatch (MMCM) test for the three datasets.

```
R> MMCM(setosa, versicolor, virginica)
```

Approximative MMCM test

```
data:  setosa, versicolor, virginica
chisq = 129.78, df = 3, p-value < 2.2e-16
alternative hypothesis: At least one pair of distributions are unequal.
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(setosa, versicolor, virginica, method = "MMCM")
```

Approximative MMCM test

```
data:  setosa, versicolor, virginica
chisq = 129.78, df = 3, p-value < 2.2e-16
alternative hypothesis: At least one pair of distributions are unequal.
```

The MMCM statistic value on its own is hard to interpret. However, the test rejects the null hypothesis of equal distributions with $p < 2.2 \cdot 10^{-16}$. Therefore, we can conclude that the observed MMCM value presents an extreme value when assuming the null. Thus the datasets are dissimilar.

3.3. Exactly two numeric datasets with target variables

The `segmentationData` dataset (Hill, LaPan, Li, and Haney 2007) in the `caret` package (Kuhn and Max 2008) includes cell body segmentation data. The dataset contains 119 imaging

measurements of 2019 cells to predict the segmentation that is divided into the two classes PS for ‘poorly segmented’ and WS for ‘well segmented’. Moreover, there is a division into 1009 observations used for training and 1010 observations used as a test set. We compare this training and test set. Ideally, the distributions of the training and test sets should be equal.

```
R> data(segmentationData, package = "caret")
R> test <- segmentationData[segmentationData$Case == "Test", -(1:2)]
R> train <- segmentationData[segmentationData$Case == "Train", -(1:2)]
```

To check the similarity of the training and test sets, we apply the method of [Ntoutsis *et al.* \(2008\)](#). For demonstration, we use all three proposed similarity measures NTO1, NTO2, and NTO3. In all cases, we do not tune the decision trees that are used to define the partitions. The `target1` and `target2` arguments have to be specified as the column names of the target variable in the first and second supplied dataset, respectively. Here, the target variable is named "Class" in both cases.

```
R> NKT(train, test, target1 = "Class", target2 = "Class", tune = FALSE)
```

Data similarity according to Ntoutsis *et al.* (2008), version 1

```
data:  train and test
s = 0.96931
alternative hypothesis: The distributions of train and test are unequal.
```

```
R> NKT(train, test, target1 = "Class", target2 = "Class", tune = FALSE,
+       version = 2)
```

Data similarity according to Ntoutsis *et al.* (2008), version 2

```
data:  train and test
s = 0.92444
alternative hypothesis: The distributions of train and test are unequal.
```

```
R> NKT(train, test, target1 = "Class", target2 = "Class", tune = FALSE,
+       version = 3)
```

Data similarity according to Ntoutsis *et al.* (2008), version 3

```
data:  train and test
s = 0.96648
alternative hypothesis: The distributions of train and test are unequal.
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(train, test, method = "NKT", target1 = "Class",
+                target2 = "Class", tune = FALSE)
```

Data similarity according to Ntoutsis et al. (2008), version 1

```
data: train and test
s = 0.96931
alternative hypothesis: The distributions of train and test are unequal.

R> DataSimilarity(train, test, method = "NKT", target1 = "Class",
+                 target2 = "Class", tune = FALSE, version = 2)
```

Data similarity according to Ntoutsis et al. (2008), version 2

```
data: train and test
s = 0.92444
alternative hypothesis: The distributions of train and test are unequal.

R> DataSimilarity(train, test, method = "NKT", target1 = "Class",
+                 target2 = "Class", tune = FALSE, version = 3)
```

Data similarity according to Ntoutsis et al. (2008), version 3

```
data: train and test
s = 0.96648
alternative hypothesis: The distributions of train and test are unequal.
```

We observe high similarity between the training and test datasets with all three methods, reflected by the similarity values s that are all close to the maximal value 1. For the method of Ntoutsis et al. (2008), no test is proposed and therefore, no p value is calculated.

3.4. Exactly two categorical datasets without target variables

The `banque` dataset from the `ade4` package (Dray and Dufour 2007) consists of bank survey data of 810 customers. All variables are categorical and contain socio-economic information of the customers. We divide the data into bank card owners and non-bank card owners and compare these two groups. In total, 243 out of the 810 customers own a bank card.

```
R> data(banque, package = "ade4")
R> card <- banque[banque$cableue == "oui", -7]
R> no.card <- banque[banque$cableue == "non", -7]
```

We use the random forest test of Hediger et al. (2022) to compare these two groups. For easier interpretation, we look at the overall out-of-bag (OOB) prediction error instead of the per-class OOB prediction error.

```
R> HMN(card, no.card, n.perm = 1000, statistic = "OverallOOB")
```

Permutation OverallOOB random forest based two-sample test

```
data: card and no.card
p.hat = 0.1605, p-value = 0.000999
alternative hypothesis: The distributions of card and no.card are unequal.
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(card, no.card, method = "HMN", n.perm = 1000,
+               statistic = "Overall100B")
```

Permutation Overall100B random forest based two-sample test

```
data:  card and no.card
p.hat = 0.1605, p-value = 0.000999
alternative hypothesis: The distributions of card and no.card are unequal.
```

The overall OOB prediction error is 0.161, which is considerably smaller than the naive prediction error of $243/810 = 0.3$. Therefore, the random forest can distinguish between the datasets, so we can conclude that the datasets differ. This is also reflected by the p value of $9.990\text{e-}04$.

3.5. More than two categorical datasets without target variables

We consider the `banque` dataset from the `ade4` package (Dray and Dufour 2007) again. This time, we split it into the nine socio-professional categories given by ‘csp’.

```
R> data(banque, package = "ade4")
R> agric <- banque[banque$csp == "agric", -1]
R> artis <- banque[banque$csp == "artis", -1]
R> cadsu <- banque[banque$csp == "cadsu", -1]
R> inter <- banque[banque$csp == "inter", -1]
R> emplo <- banque[banque$csp == "emplo", -1]
R> ouvri <- banque[banque$csp == "ouvri", -1]
R> retra <- banque[banque$csp == "retra", -1]
R> inact <- banque[banque$csp == "inact", -1]
R> etudi <- banque[banque$csp == "etudi", -1]
```

We apply the classifier two-sample test (C2ST). First, we use the default K -NN classifier. Categorical variables are dummy-coded.

```
R> C2ST(agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi)
```

Approximative Classifier Two-Sample Test using knn

```
data:  agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi
p.hat = 0.33333, size = 567.00000, prob = 0.22593, p-value =
1.078e-07
alternative hypothesis: At least one pair of distributions are unequal.
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(agric, artis, cadsu, inter, emplo, ouvri, retra, inact,
+               etudi, method = "C2ST")
```

Approximative Classifier Two-Sample Test using knn

```
data:  agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi
p.hat = 0.33333, size = 567.00000, prob = 0.22593, p-value =
1.078e-07
alternative hypothesis: At least one pair of distributions are unequal.
```

The accuracy of the K -NN classifier is 0.333. It is larger than the naive accuracy for always predicting the largest class, which is given by `prob = 0.226` in the output. The classifier seems to be able to distinguish between the datasets, and we can therefore regard them as dissimilar. Moreover, the null hypothesis of equal distributions can be rejected with a p value of `1.078e-07`.

For demonstration, we additionally perform the C2ST with a neural net classifier.

```
R> C2ST(agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi,
+       classifier = "nnet", train.args = list(trace = FALSE))
```

Approximative Classifier Two-Sample Test using nnet

```
data:  agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi
p.hat = 0.27778, size = 567.00000, prob = 0.22593, p-value =
2.425e-05
alternative hypothesis: At least one pair of distributions are unequal.
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(agric, artis, cadsu, inter, emplo, ouvri, retra, inact,
+                etudi, method = "C2ST", classifier = "nnet",
+                train.args = list(trace = FALSE))
```

Approximative Classifier Two-Sample Test using nnet

```
data:  agric, artis, cadsu, inter, emplo, ouvri, retra, inact, etudi
p.hat = 0.30556, size = 567.00000, prob = 0.22593, p-value =
1.826e-06
alternative hypothesis: At least one pair of distributions are unequal.
```

The results are very similar to using K -NN.

3.6. Exactly two categorical datasets with target variables

We consider the `banque` dataset from the `ade4` package (Dray and Dufour 2007) again. In this case, we interpret the savings bank amount (`eparliv`) variable as the target variable, which is again supplied via the `target1` and `target2` arguments. It is divided into the three categories ‘> 20000’, ‘> 0 and < 20000’, and ‘nulle’. We divide the data into the socio-professional categories as before. We use the optimal transport dataset distance (OTDD) to compare the resulting datasets for craftsmen, shopkeepers, company directors (‘`artis`’), to

that of higher intellectual professions ('cadsu') and to that of manual workers ('ouvri'). As all variables are categorical, we use the Hamming distance instead of the default Euclidean distance.

```
R> OTDD(artis, cadsu, target1 = "eparliv", target2 = "eparliv",
+       feature.cost = hammingDist)
```

Optimal Transport Dataset Distance

```
data:  artis and cadsu
OTDD = 44.166
alternative hypothesis: Distributions of artis and cadsu are unequal
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(artis, cadsu, method = "OTDD", target1 = "eparliv",
+               target2 = "eparliv", feature.cost = hammingDist)
```

Optimal Transport Dataset Distance

```
data:  artis and cadsu
OTDD = 44.166
alternative hypothesis: Distributions of artis and cadsu are unequal
```

We obtain a dataset distance of 44.166 between craftsmen/shopkeepers/company directors and executives/higher intellectual professions. For the OTDD, low values correspond to high similarity, and the minimum value is 0. The observed value is clearly larger than zero, so the datasets are not exactly similar. How dissimilar they are is however hard to interpret from the observed OTDD value on its own. For the OTDD, no test is proposed and therefore, no p value is calculated.

```
R> OTDD(artis, ouvri, target1 = "eparliv", target2 = "eparliv",
+       feature.cost = hammingDist)
```

Optimal Transport Dataset Distance

```
data:  artis and ouvri
OTDD = 49.427
alternative hypothesis: Distributions of artis and ouvri are unequal
```

Alternatively, we can always use the `DataSimilarity()` function:

```
R> DataSimilarity(artis, ouvri, method = "OTDD", target1 = "eparliv",
+               target2 = "eparliv", feature.cost = hammingDist)
```

Optimal Transport Dataset Distance

```
data:  artis and ouvri
OTDD = 49.427
alternative hypothesis: Distributions of artis and ouvri are unequal
```


We obtain a dataset distance of 49.427 between craftsmen/shopkeepers/company directors and manual workers. Again, this value on its own is hard to interpret. However, we can compare the values and conclude that the data of craftsmen/shopkeepers/company directors is more similar to that of executives/higher intellectual professions than to that of manual workers.

4. Implementation overview

Table 2 gives an overview of all wrapper functions included in the package. For each method, the original implementation, the new function name, and the applicability to data with a target variable, numerical data, categorical data, and multiple samples are given. Note that the applicability statements refer to the specific implementation of the method. Some of the methods are in theory applicable to a broader range of data types than are implemented. Moreover, note that most implementations are only applicable to either numerical or categorical data except for the classifier-based methods `HMN()` and `C2ST()`, which can handle both data types simultaneously as long as the selected classifier can do so. The `MMD()` implementation can also handle both data types, but a matching kernel function has to be implemented. Note that the graph-based tests cannot deal with both numerical and categorical data due to ties, even if a distance function that can handle both is supplied. More details on the methods and their implementation can be found in Section 5.

Table 3 gives an overview of the newly implemented methods and their applicability. A few of these methods were already implemented in another programming language, as described in the implementation details in Section 5.

Method	Original function	New function	y	Num	Cat	$k > 2$
KMD (Huang and Sen 2023)	KMD::KMD(), KMD::KMD_test() (Huang 2022)	KMD()	✗	✓	✗*	✓
Friedman and Rafsky (1979)	gTests::g.tests() (Chen and Zhang 2017)	FR()	✗	✓	✓	✗
Cross-match test (Rosenbaum 2005)	crossmatch::crossmatch() (Heller, Small, and Rosen- baum 2024)	Rosenbaum()	✗	✓	✗	✗
Cramér test (Baring- haus and Franz 2004)	cramer::cramer.test() (Franz 2024)	Cramer()	✗	✓	✗	✗
Energy statistic (Székely and Rizzo 2017)	energy::eqdist.test() (Rizzo and Szekely 2024)	Energy()	✗	✓	✗	✓
Hediger <i>et al.</i> (2022)	hypoRF::hypoRF() (Hediger, Michel, and Naef 2024)	HMN()	✗	✓	✓	✗
Baringhaus and Franz (2010)	cramer::cramer.test() (Franz 2024)	BF()	✗	✓	✗	✗
Bahr (1996)	cramer::cramer.test() (Franz 2024)	Bahr()	✗	✓	✗	✗
Wasserstein distance	Ecume::wasserstein_permut() (Roux de Bezieux 2024)	Wasserstein()	✗	✓	✗	✗
Chen and Friedman (2017)	gTests::g.tests() (Chen and Zhang 2017)	CF()	✗	✓	✓	✗
Chen, Chen, and Su (2018)	gTests::g.tests() (Chen and Zhang 2017)	CCS()	✗	✓	✓	✗
Ball divergence (Pan, Tian, Wang, and Zhang 2018)	Ball::bd.test() (Zhu, Pan, Zheng, and Wang 2021)	BallDivergence()	✗	✓	✗	✓
Song and Chen (2022)	gTestsMulti::gtestsmulti() (Song and Chen 2023b)	SC()	✗	✓	✗	✓
DISCO (Rizzo and Székely 2010)	energy::eqdist.test() (Rizzo and Szekely 2024)	DISCOB(), DISCOF()	✗	✓	✗	✓
Zhang and Chen (2019)	gTests::g.tests() (Chen and Zhang 2017)	ZC()	✗	✓	✓	✗
RI test (Paul, De, and Ghosh 2022b)	HDLSSkST::RIttest() (Paul, De, and Ghosh 2022a)	RIttest()	✗	✓	✗	✓
FS test (Paul <i>et al.</i> 2022b)	HDLSSkST::FSttest() (Paul <i>et al.</i> 2022a)	FSttest()	✗	✓	✗	✓
Maximum Mean Dis- crepancy (MMD) (Gret- ton, Borgwardt, Rasch, Schölkopf, and Smola 2006)	kernlab::kmmd() (Karat- zoglou, Smola, Hornik, and Zeileis 2004)	MMD()	✗	✓	✗*	✗
Song and Chen (2023a)	kerTests::kertests() (Song and Chen 2023c)	GPk()	✗	✓	✗*	✗
Mukhopadhyay and Wang (2020b)	LPKsample::GLP() (Mukhopadhyay and Wang 2020a)	MW()	✗	✓	✗*	✓
Chen, Dou, and Qiao (2013)	gTests::g.tests_cat() (Chen and Zhang 2017)	FR_cat(), CF_cat(), CCS_cat(), ZC_cat()	✗	✓	✓	✗

Classifier	Two-Sample	Ecume::classifier_test()	C2ST()	✗	✓	✓	✓
Test	(Lopez-Paz and Oquab 2017)	(Roux de Bezieux 2024)					

Table 2: Implemented wrapper functions. y : Can the method deal with a target variable in the dataset? Num: Is the method as implemented applicable to numeric data? Cat: Is the method as implemented applicable to categorical data? $k > 2$: Is the method as implemented applicable to more than two datasets at a time? ✗*: Method is, in theory, applicable, but implementation is not. ✓*: Implementation is applicable although this case is not described in the literature.

Method	New function	y	Num	Cat	$k > 2$
Mukherjee <i>et al.</i> (2022)	MMCM()	✗	✓	✓*	✓
Petrie (2016)	Petrie()	✗	✓	✓*	✓
Biswas, Mukhopadhyay, and Ghosh (2014)	BMG()	✗	✓	✗	✓
Deb and Sen (2021)	DS()	✗	✓	✗	✗
Ntoutsis <i>et al.</i> (2008)	NKT()	✓	✓	✗	✗
Ganti, Gehrke, Ramakrishnan, and Loh (1999)	GGRL()	✓	✓	✗*	✗
Alvarez-Melis and Fusi (2020)	OTDD()	✓	✓	✓	✗
Jeffreys divergence	Jeffreys()	✗	✓	✗	✗
Biswas and Ghosh (2014)	BG2()	✗	✓	✗	✗
Engineer metric	engineerMetric()	✗	✓	✗	✗
Schilling (1986) and Henze (1988)	SH()	✗	✓	✗	✗
Barakat, Quade, and Salama (1996)	BQS()	✗	✓	✗	✗
Yu, Martin, Rothman, Zheng, and Lan (2007)	YMRZL()	✗	✓	✓	✗
Li, Hu, and Zhang (2022)	LHZ()	✗	✓	✗	✗
Constrained Minimum Distance (Tatti 2007)	CMDistance()	✗	✗	✓	✗
Biau and Györfi (2005)	BG()	✗	✓	✗	✓
DiProPerm test (Wei <i>et al.</i> 2016)	DiProPerm()	✗	✓	✗	✗

Table 3: Newly implemented functions. y : Can the method deal with a target variable in the dataset? Num: Is the method as implemented applicable to numeric data? Cat: Is the method as implemented applicable to categorical data? $k > 2$: Is the method as implemented applicable to more than two datasets at a time? ✗*: Method is, in theory, applicable, but implementation is not. ✓*: Implementation is applicable although this case is not described in the literature.

5. Implementation details

5.1. KMD (Huang and Sen 2023)

The *kernel measure of multi-sample dissimilarity* (KMD) introduced by Huang and Sen (2023) is a kernel-based test using the association between the variables and the sample membership to quantify the dissimilarity of multiple samples. Denote the dataset membership of each point in the pooled sample $\{Z_1, \dots, Z_N\}$ by $\{\Delta_1, \dots, \Delta_N\}$. $\{(\Delta_i, Z_i)\}_{i=1}^N$ can be seen as an i.i.d. sample from $(\tilde{\Delta}, \tilde{Z})$ with distribution μ defined by $P(\tilde{\Delta} = i) = \pi_i, i = 1, \dots, k$ and $\tilde{Z}|\tilde{\Delta} = i \sim F_i$. Let $(\tilde{Z}_1, \tilde{\Delta}_1), (\tilde{Z}_2, \tilde{\Delta}_2)$ be i.i.d. samples from μ and $(\tilde{Z}, \tilde{\Delta}), (\tilde{Z}, \tilde{\Delta}') \sim \mu$ with $\tilde{\Delta}, \tilde{\Delta}'$ conditionally independent given \tilde{Z} . Denote by K a kernel function over $\{1, \dots, k\}$, e.g., the discrete kernel $K(x, y) := \mathbb{1}(x = y)$. Then the KMD is defined as

$$\eta(P_1, \dots, P_k) := \frac{\mathbb{E}[K(\tilde{\Delta}, \tilde{\Delta}')] - \mathbb{E}[K(\tilde{\Delta}_1, \tilde{\Delta}_2)]}{\mathbb{E}[K(\tilde{\Delta}, \tilde{\Delta})] - \mathbb{E}[K(\tilde{\Delta}_1, \tilde{\Delta}_2)]}.$$

It can be estimated using a similarity graph \mathcal{G} , e.g., the K -nearest neighbor graph or the minimum spanning tree (MST), on the pooled sample. Denote by $(Z_i, Z_j) \in \mathcal{E}(\mathcal{G})$ that there is an edge in \mathcal{G} connecting Z_i and Z_j . Moreover, let o_i be the out-degree of Z_i in \mathcal{G} . Then an estimator for η is defined as

$$\hat{\eta} := \frac{\frac{1}{N} \sum_{i=1}^N \frac{1}{o_i} \sum_{j:(Z_i, Z_j) \in \mathcal{E}(\mathcal{G})} K(\Delta_i, \Delta_j) - \frac{1}{N(N-1)} \sum_{i \neq j} K(\Delta_i, \Delta_j)}{\frac{1}{N} \sum_{i=1}^N K(\Delta_i, \Delta_i) - \frac{1}{N(N-1)} \sum_{i \neq j} K(\Delta_i, \Delta_j)}.$$

An asymptotic and a permutation k -sample test are proposed based on the KMD.

The implementation of the new function `KMD()` combines the calculation of KMD and the corresponding p value using the functions `KMD()` and `KMD_test()`, respectively, from the **KMD** package (Huang 2022). Moreover, the inputs of the new function are simply the individual datasets instead of the pooled data matrix and sample IDs. By default, the asymptotic test is performed (`n.perm = 0`) using the K -nearest neighbor graph with $K = \lceil N/10 \rceil$, where N denotes the total sample size of the pooled sample, and a discrete kernel. The options for the graph are restricted to `knn` and `mst` by the implementations from the **KMD** package. A user-specified kernel can be used only when a kernel matrix is supplied instead of the keyword “discrete” for the `kernel` argument of the new function.

5.2. Edge-count tests (Friedman and Rafsky 1979; Chen and Zhang 2013; Chen *et al.* 2018; Chu and Chen 2019)

The tests by Friedman and Rafsky (1979), Chen and Friedman (2017), Chen *et al.* (2018), and Chu and Chen (2019) are graph-based two-sample tests that use the edge counts in a similarity graph like the (K) -MST on the pooled sample. They make use of the number of edges that connect points within the first sample, R_1 , the number of edges that connect points within the second sample, R_2 , and the number of edges that connect points from different samples R_{12} . The original edge-count test by Friedman and Rafsky (1979) takes the standardized between-sample edge-count

$$T_{\text{FR}} = \frac{R_{12} - \mathbb{E}_{H_0}(R_{12})}{\sqrt{\text{VAR}_{H_0}(R_{12})}}$$

as its test statistic. The expectation and variance under the null can be calculated analytically. Chen and Friedman (2017) noted that this has low power against scale alternatives and proposed the *generalized edge-count test* using

$$T_{\text{CF}} = (R_1 - \mathbb{E}_{H_0}(R_1), R_2 - \mathbb{E}_{H_0}(R_2)) \text{COV}_{H_0}^{-1} \left(\begin{pmatrix} R_1 \\ R_2 \end{pmatrix} \right) \begin{pmatrix} R_1 - \mathbb{E}_{H_0}(R_1) \\ R_2 - \mathbb{E}_{H_0}(R_2) \end{pmatrix}.$$

Chen *et al.* (2018) found some problems with the original edge-count test for unequal sample sizes of the two datasets, based on which they proposed the *weighted edge-count test* using the weighted edge-counts

$$R_w = \frac{n_1}{N} R_1 + \frac{n_2}{N} R_2,$$

where n_1 denotes the sample size of the first dataset and n_2 the sample size of the second dataset, and $N = n_1 + n_2$ the total sample size in the pooled sample. The *weighted edge-count test* statistic is then defined as the standardized weighted edge count

$$T_{\text{CCS}} = \frac{R_w - E_{H_0}(R_w)}{\sqrt{\text{VAR}_{H_0}(R_w)}}.$$

Lastly, the *max-type edge count* (Chu and Chen 2019) test additionally uses the difference of the edge counts in the samples, i.e.,

$$R_d = R_1 - R_2.$$

Its test statistic is defined as

$$T_{\text{ZC}} = \max \left(\kappa \frac{R_w - E_{H_0}(R_w)}{\sqrt{\text{VAR}_{H_0}(R_w)}}, \left| \frac{R_d - E_{H_0}(R_d)}{\sqrt{\text{VAR}_{H_0}(R_d)}} \right| \right),$$

where κ is a constant that has to be chosen prior to performing the test. $\kappa \in \{1.31, 1.14, 1\}$ is recommended based on a small power simulation for normal data with shift or scale alternatives.

Wrapper functions around `g.tests()` from the `gTests` package (Chen and Zhang 2017) are implemented. These do not need a pre-calculated graph as input but allow specifying a distance function (`dist.fun`) and a function for calculating a similarity graph (`graph.fun`) and then calculating the similarity graph internally. The new input also includes both datasets. We find this more intuitive and less error-prone than supplying an edge matrix and two vectors of indices specifying the dataset membership as for the original `g.tests()` function. The new implementation forces the user to choose one of the tests first and then perform it, instead of performing all tests at once. Moreover, the users have to decide whether they want to perform the permutation test or the approximative test.

For the Friedman-Rafsky test, there is an additional implementation in the `GSAR` package, but there the test statistic is standardized by the empirical mean and standard deviation rather than the theoretical mean and standard deviation of the test statistic under the null hypothesis as proposed in the original article. Therefore, we use the `gTests` implementation here.

5.3. Edge-count tests for categorical data (Chen and Zhang 2013; Zhang and Chen 2019)

These methods are adaptations of the previously mentioned edge-count tests for categorical data. With categorical data, the problem of ties in the distance matrix arises. Ties lead to non-unique solutions for the similarity graph construction and therefore also to non-unique values of the proposed test statistics. This can be solved by either taking the union of all optimal graphs and calculating the respective statistic on this union graph or by averaging the test statistic values over all optimal graphs. The new implementation of the categorical graph-based tests is again a wrapper function that includes the calculation of the edge matrix. For this, the function `getGraph()` from the `gTests` package is used. Therefore, the choice of the similarity graph is restricted to the K -nearest neighbors and the K -MST. Still, a distance function can be supplied. By default, this is the sum of unequal classes. The calculation of the frequency table of all observations and the similarity graph on this are performed internally; thus, again, only the datasets have to be supplied by the user. Moreover, the method for aggregating the graphs has to be supplied. Possible options are averaging ("`a`") and union ("`u`") over graphs.

5.4. Cross-match test (Rosenbaum 2005)

The Rosenbaum cross-match test uses a similar approach as the Friedman-Rafsky test, but based on the optimal non-bipartite matching instead of the MST as a similarity graph (see Section 1.3). The new function `Rosenbaum()` is a wrapper around the `crossmatchtest()` function from the `crossmatch` package (Heller et al. 2024). Again, a distance function can be supplied. By default, this is `stats::dist()`, i.e., the Euclidean distance. The new function then calculates the distance matrix internally. Again, we find this more straightforward from a user perspective than supplying a distance matrix on the pooled sample and a vector specifying the dataset membership of each observation. The output of the function includes the raw edge count, its standard error, and expectation under the null like for the `crossmatch` implementation. In contrast, only either the exact or the approximative p value is returned. By default (`exact = TRUE`), the exact p value is returned.

This is appropriate for samples that are not too large. Note that with a pooled sample size of 340 or more, it is numerically impossible to derive the exact distribution due to the factorials involved in the calculation, and `crossmatchtest()` will return a missing value for the exact p value.

5.5. Energy statistic and generalizations by Baringhaus and Franz (2010)

The energy statistic is a popular two- and k -sample statistic based on interpoint distances. The k -sample statistic is defined as

$$T_{\text{Energy}} = \sum_{1 \leq i < j \leq k} \frac{n_i n_j}{n_i + n_j} \left(\frac{2}{n_i n_j} \sum_{u=1}^{n_i} \sum_{v=1}^{n_j} \|X_u^{(i)} - X_v^{(j)}\| \right. \\ \left. - \frac{1}{n_i^2} \sum_{u=1}^{n_i} \sum_{v=1}^{n_i} \|X_u^{(i)} - X_v^{(i)}\|_2 - \frac{1}{n_j^2} \sum_{u=1}^{n_j} \sum_{v=1}^{n_j} \|X_u^{(j)} - X_v^{(j)}\|_2 \right).$$

For a comprehensive review of the literature on the energy statistic and its applications, please refer to Székely and Rizzo (2017). A permutation test can be performed based on the energy statistic. In the two-sample case, the energy statistic is equal to two times the Cramér test statistic of Baringhaus and Franz (2004) and therefore the tests are equivalent. However, a Bootstrap instead of a permutation test is proposed for the Cramér test. Baringhaus and Franz (2010) propose a test statistic that generalizes the Cramér test statistic by using a continuous function ϕ such that $\phi(\|x - y\|^2)$ is a negative definite kernel instead of the Euclidean distances. Different examples for ϕ are given, including as special cases the Cramér test, the test by Bahr (1996), and the test by Szabo, Boucher, Carroll, Klebanov, Tsodikov, and Yakovlev (2002). Overall, $\phi(z) = \log(1 + z)$ is recommended for general alternatives based on a simulation study, and the Cramér test is recommended for location alternatives. The tests of Baringhaus and Franz (2010) are implemented in the **cramer** package (Franz 2024). The new implementation is a simple wrapper to unify input and output naming and types. The energy statistic is implemented in the **energy** package (Rizzo and Szekely 2024). For the corresponding wrapper, the input type was changed more since the original implementation had the pooled sample and the sample sizes as the input. The **energy** implementation outsourced the calculation of the energy statistic to **C**, which gives it a notable advantage with regard to computing time over the **cramer** implementation.

5.6. Random forest-based test (Hediger et al. 2022)

The random forest-based method of Hediger et al. (2022) is briefly described above in Section 1.3. The function here is a wrapper around the `hypoRF()` function from the **hypoRF** package (Hediger et al. 2024) that only renames arguments for consistency with the other methods. Note that the implemented per-class OOB statistics differ for the permutation test and the approximate test: for the permutation test, the sum of the per-class OOB errors is returned, for the asymptotic version, the standardized sum is returned.

5.7. Wasserstein distance

The q -Wasserstein distance (Vaserstein 1969) of two distributions F_1 and F_2 on \mathcal{X} is defined as

$$W(F_1, F_2) := \left(\min_{\pi \in \Pi(F_1, F_2)} \int_{\mathcal{X} \times \mathcal{X}} d_{\mathcal{X}}(x, y)^q d\pi(x, y) \right)^{1/q},$$

where $d_{\mathcal{X}}$ is the metric that \mathcal{X} is provided with, and

$$\Pi(F_1, F_2) := \{\pi_{1,2} \in \mathcal{P}(\mathcal{X} \times \mathcal{X}) | \pi_1 = F_1, \pi_2 = F_2\}$$

is the set of joint distributions over the product space $\mathcal{X} \times \mathcal{X}$ with marginal distributions F_1 and F_2 .

In the **Ecume** package (Roux de Bezieux 2024), a permutation test based on the Wasserstein distance is implemented.

5.8. Ball divergence (Pan et al. 2018)

The Ball divergence measures the difference between two probability measures. It is defined as the square of the measure difference over a given closed ball collection. It can be estimated as

$$\widehat{\text{BD}} = A + C,$$

where

$$A = \frac{1}{n_1^2} \sum_{i,j=1}^{n_1} \left(A_{ij}^{(1)} - A_{ij}^{(2)} \right)^2,$$

$$C = \frac{1}{n_2^2} \sum_{l,m=1}^{n_2} \left(C_{lm}^{(1)} - C_{lm}^{(2)} \right)^2,$$

and

$$A_{ij}^{(1)} = \frac{1}{n_1} \sum_{u=1}^{n_1} \mathbb{1}(X_u^{(1)} \in \bar{B}(X_i^{(1)}, d(X_i^{(1)}, X_j^{(1)}))),$$

$$A_{ij}^{(2)} = \frac{1}{n_2} \sum_{v=1}^{n_2} \mathbb{1}(X_v^{(2)} \in \bar{B}(X_i^{(1)}, d(X_i^{(1)}, X_j^{(1)}))),$$

$$C_{lm}^{(1)} = \frac{1}{n_1} \sum_{u=1}^{n_1} \mathbb{1}(X_u^{(1)} \in \bar{B}(X_l^{(2)}, d(X_l^{(2)}, X_m^{(2)}))),$$

$$C_{lm}^{(2)} = \frac{1}{n_2} \sum_{v=1}^{n_2} \mathbb{1}(X_v^{(2)} \in \bar{B}(X_l^{(2)}, d(X_l^{(2)}, X_m^{(2)}))),$$

with $\bar{B}(X_i^{(l)}, d(X_i^{(l)}, X_j^{(l)}))$ denoting the closed Ball around $X_i^{(l)}$ with radius equal to the distance d of the points $X_i^{(l)}$ and $X_j^{(l)}$, $l \in \{1, 2\}$. Therefore, the first part of the Ball divergence, A , consists of squared distances of proportions of data points from the first sample lying within closed balls around data points from the first sample and of data points from the second sample lying within closed balls around data points from the first sample. The second part, C , consists of squared distances of proportions of data points from the first sample lying within closed balls around data points from the second sample and of data points from the second sample lying within closed balls around data points from the second sample. For both parts, the mean over all such Balls with radii equal to the distances of the center point of the ball to all other points from the same sample is taken. For multiple samples, the pairwise test statistics can be summarized by summing up the pairwise divergences, or by taking the maximum of sums of the Ball divergences from each sample to all other samples, or by summing the largest $k - 1$ pairwise Ball divergences.

The implementation here is a wrapper around the `bd.test()` function from the **Ball** package (Zhu *et al.* 2021). In contrast to the original implementation, the new wrapper returns an object of class ‘**htest**’ in the multi-sample case, although in that case no test is conducted. Moreover, only the summarized statistic according to the specified `kbd.type`, which determines how the pairwise Ball divergences are summarized, is returned.

5.9. Multisample graph-based tests (Song and Chen 2022)

Song and Chen (2022) propose three new tests for the k -sample problem that use the between-sample edges and the within-sample edges of a similarity graph on the pooled sample. Let R^W denote the vector containing the numbers of within-sample edges for each of the k samples and R^B denote the vector containing the numbers of between-sample edges for all $k(k - 1)$ pairs of different samples. Then the first test statistic is given by

$$S = S^W + S^B$$

$$S^W = \left(R^W - \mathbb{E}_{H_0}(R^W) \right)^\top \text{COV}_{H_0}^{-1} \left(R^W \right) \left(R^W - \mathbb{E}_{H_0}(R^W) \right)$$

$$S^B = \left(R^B - \mathbb{E}_{H_0}(R^B) \right)^\top \text{COV}_{H_0}^{-1} \left(R^B \right) \left(R^B - \mathbb{E}_{H_0}(R^B) \right).$$

The second test statistic is based on the vector R^A of all linearly independent numbers of edges between and within samples, i.e., all numbers of edges between all pairs of samples, including the pairs of a sample with itself, except for the pair of sample $(k - 1)$ and sample k . The test statistic is then defined as

$$S^A = \left(R^A - \mathbb{E}_{H_0}(R^A) \right)^\top \text{COV}_{H_0}^{-1} \left(R^A \right) \left(R^A - \mathbb{E}_{H_0}(R^A) \right).$$

All expectations and covariances under the null can be calculated analytically again. While $\text{COV}_{H_0}(R^W)$ is shown to be always invertible, no such proof exists for $\text{COV}_{H_0}(R^B)$ and $\text{COV}_{H_0}(R^A)$. Therefore, Song and

Chen (2022) suggest checking the invertability numerically before applying the test and using a generalized inverse if necessary. This is already done within their implementation. Based on S^A , an asymptotic test can easily be performed. The asymptotic distribution of S is more complicated and hard to compute in practice, therefore, a fast test is suggested instead. It combines the tests using S^W and S^B and takes the Bonferroni-adjusted p value of both these tests. Alternatively, a permutation test can be performed for either S^A or S . The implementation here for the test of Song and Chen (2022) is a wrapper around the `gtestsmulti()` function from `gTestsMulti` (Song and Chen 2023b). The input is simplified as for the wrapper around `g.tests()`. The user has to choose whether the original (S) or the fast (S^A) version of the test should be performed. If the number of permutations for the permutation test (`n.perm`) is set to 0, the approximate test is performed; otherwise, the permutation p value is reported.

5.10. DISCO

Rizzo and Székely (2010) show that the energy test can be seen as the treatment sum of squares in an ANOVA interpretation of the k -sample problem. As the measure of dispersion for univariate or multivariate responses based on all pairwise distances between-sample elements for ANOVA

$$d_\alpha(X^{(1)}, X^{(2)}) = \frac{n_1 n_2}{n_1 + n_2} [2g_\alpha(X^{(1)}, X^{(2)}) - g_\alpha(X^{(1)}, X^{(1)}) - g_\alpha(X^{(2)}, X^{(2)})]$$

is proposed with

$$g_\alpha(X^{(1)}, X^{(2)}) = \frac{1}{n_1 n_2} \sum_{u=1}^{n_1} \sum_{v=1}^{n_2} \|X_u^{(1)} - X_v^{(2)}\|_2^\alpha.$$

With this, Rizzo and Székely (2010) derive their so-called *distance components (DISCO) decomposition* for $\alpha \in (0, 2]$. It partitions the total dispersion in the samples

$$T_\alpha = \frac{N}{2} g_\alpha(Z, Z),$$

into components

$$T_\alpha = S_\alpha + W_\alpha$$

analogous to the variance components in ANOVA. Here, Z denotes the pooled sample and the between-sample measure of dispersion S_α and the within-sample measure of dispersion W_α , respectively, are defined as

$$S_\alpha = \sum_{1 \leq i < j \leq k} \frac{n_i + n_j}{2N} d_\alpha(X^{(i)}, X^{(j)}),$$

$$W_\alpha = \sum_{i=1}^k \frac{n_i}{2} g_\alpha(X^{(i)}, X^{(i)}).$$

The between-sample measure of dispersion S_α can be used directly to compare the distributions in a k -sample permutation test (`DISCOB()`). Alternatively, the statistic

$$F_\alpha = \frac{S_\alpha / (k - 1)}{W_\alpha / (N - k)}$$

can be used in a k -sample permutation test (`DISCOF()`). For each index $\alpha \in (0, 2)$, this determines a nonparametric test for the multi-sample problem that is statistically consistent against general alternatives. For $\alpha = 2$, it equals the usual ANOVA F -test. The choice of the index α is difficult. In general, the computational costs for calculating Gini means g_α , in terms of which the test statistic can be formulated, are $\mathcal{O}(N^2)$. For $\alpha = 1$, it can be linearized and computation time reduces to $\mathcal{O}(N \log N)$. The simplest and most natural choice for α is one. For heavy-tailed distributions, a small α is recommended.

The test is implemented by permutation Bootstrap in the R package `energy` (Rizzo and Székely 2024). The new implementations of the between-sample and of the DISCO F -test are wrappers that mainly unify the inputs and outputs, which differed between the two tests in the original implementation. Moreover, the input format is again changed from the pooled sample and the dataset labels to the individual datasets.

5.11. (Modified / multiscale / aggregated) RI and FS test

Paul *et al.* (2022b) propose distribution-free k -sample tests intended for the high dimension low sample size (HDLSS) setting. The tests are based on clustering the pooled sample and comparing the resulting clustering

to the true dataset membership via a contingency table. If the datasets come from the same distribution, the cluster and dataset membership are independent, while if the datasets come from different distributions, the clustering depends on the true dataset membership. As a clustering algorithm, Paul *et al.* (2022b) suggest using K -means based on the generalized version of the *mean absolute difference of distances* (MADD)

$$\rho_{h,\varphi}(z_i, z_j) = \frac{1}{N-2} \sum_{m \in \{1, \dots, N\} \setminus \{i, j\}} |\varphi_{h,\psi}(z_i, z_m) - \varphi_{h,\psi}(z_j, z_m)|,$$

as proposed by Sarkar and Ghosh (2020) for the HDLSS setting. Here, $z_i, i = 1, \dots, N$, denote realizations from the pooled sample and

$$\varphi_{h,\psi}(z_i, z_j) = h \left(\frac{1}{p} \sum_{l=1}^p \psi |z_{il} - z_{jl}| \right),$$

where $h : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ and $\psi : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ are continuous and strictly increasing functions. $\psi(t) = t^2$, $\psi(t) = 1 - \exp(-t)$, $\psi(t) = 1 - \exp(-t^2)$, $\psi(t) = \log(1 + t)$, and $\psi(t) = t$ are considered in combination with $h(t) = \sqrt{t}$ and $h(t) = t$. The number of clusters has to be chosen in advance. A natural choice is to set the number of clusters to k . For the RI test, the Rand index of the clustering is used as a test statistic. It is zero when the clustering is perfect, i.e., when the cluster membership is a permutation of the true dataset membership. The test rejects for low values since the Rand index should take higher values when all clusters have similar distributions of class labels. The critical value can be calculated using a generalized hypergeometric distribution. Due to the discreteness of the Rand index, Paul *et al.* (2022b) propose to use a randomized test. For the FS test, the generalized Fisher's test statistic for testing for independence in an $k \times \ell$ contingency table is used. Again, a randomized test using the generalized hypergeometric distribution to find the critical values is proposed.

Paul *et al.* (2022b) additionally propose modified versions of the tests (MRI, MFS test). For these, the number of clusters is estimated from the data using the Dunn index since setting the number of clusters to k might fail in case of multimodal distributions, where a larger number of clusters might be required, and then multiple clusters can correspond to one dataset.

Moreover, multiscale versions of the tests are presented (MSRI, MSFS test) for the case where the number of clusters is unclear. The respective tests are then performed for different numbers of clusters, and the results are aggregated using a Bonferroni adjustment for the individual tests. Still, an upper limit for the number of clusters to be considered must be chosen. The implementation also includes aggregated tests (AFS / ARI test) that perform all pairwise FS / MFS or RI / MRI tests, respectively, on the samples and aggregate the results by taking the minimum test statistic value and applying a multiple testing procedure.

The tests are implemented in the R package **HDLSskST** (Paul *et al.* 2022a). The main difference between the new wrapper functions and the original implementation is that the modified and multiscale versions of the RI and FS tests can be performed with the same function as the original tests. The test can be chosen via the newly introduced **version** argument of the **FStest()** and **RItest()** functions. One advantage of this is that the input and output formats are unified between the versions of the test. In the original implementation of the test, the elements of the output list differ both content-wise and in their names between the tests. Moreover, the input of the tests differs slightly between the original functions for the different tests. The input is also unified to match the input of the other functions in the **DataSimilarity** package and therefore consists simply of the datasets instead of a pooled data matrix, a vector with the dataset affiliation of each observation, and a vector of the sample sizes. We think this is easier to understand and less error-prone from a user perspective.

5.12. MMD

The *maximum mean discrepancy* (MMD) uses a kernel mean embedding to define a metric for probability distributions. Kernel mean embeddings extend feature maps ϕ to the space of probability distributions by representing each distribution F as a mean function

$$\phi(F)(\cdot) = \mu_F(\cdot) := \int_{\mathcal{X}} K(x, \cdot) dF(x) = \mathbb{E}_F(K(X, \cdot)),$$

where $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ is a symmetric and positive definite kernel function. A reproducing kernel Hilbert space (RKHS) \mathcal{H} of functions on the domain \mathcal{X} with kernel K is a Hilbert space of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ with dot product $\langle \cdot, \cdot \rangle$ that satisfies the reproducing property

$$\langle f(\cdot), K(x, \cdot) \rangle = f(x) \Rightarrow \langle K(x, \cdot), K(x', \cdot) \rangle = K(x, x'),$$

such that the linear map from a function to its value at x can be seen as an inner product. Then the kernel mean embedding as given above is a transformation of the distribution F to an element in the reproducing kernel Hilbert space (RKHS) \mathcal{H} corresponding to the kernel K (Muandet, Fukumizu, Sriperumbudur, and Schölkopf 2017). For characteristic kernels, the kernel mean representation captures all information about the distribution F , which implies $\|\mu_{F_1} - \mu_{F_2}\|_{\mathcal{H}} = 0 \Leftrightarrow F_1 = F_2$ (Fukumizu, Bach, and Jordan 2004; Sriperumbudur, Gretton, Fukumizu, Lanckriet, and Schölkopf 2008; Sriperumbudur, Gretton, Fukumizu, Schölkopf, and Lanckriet 2010). Therefore, the MMD measures the difference between two distributions as

$$\text{MMD}(\mathcal{H}, F_1, F_2) = \|\mu_{F_1} - \mu_{F_2}\|_{\mathcal{H}}.$$

Here, the implementation `kmmd()` from the `kernelab` package (Karatzoglou *et al.* 2004) is used. The alternative implementation from the `Ecume` does not include an automatic choice of the kernel parameter. The new implementation adds a permutation test to the `kernelab` implementation.

5.13. GPK (Song and Chen 2023a)

Song and Chen (2023a) propose another kernel-based test for which they decompose the squared MMD estimator as

$$\widehat{\text{MMD}}^2 = \alpha + \beta - 2\gamma,$$

where

$$\begin{aligned}\alpha &= \frac{1}{n_1(n_1 - 1)} \sum_{i=1}^{n_1} \sum_{\substack{j=1 \\ j \neq i}}^{n_1} K(X_i, X_j), \\ \beta &= \frac{1}{n_2(n_2 - 1)} \sum_{i=1}^{n_2} \sum_{\substack{j=1 \\ j \neq i}}^{n_2} K(Y_i, Y_j), \\ \gamma &= \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} K(X_i, Y_j).\end{aligned}$$

As a new statistic, they propose to use

$$\text{GPK} = (\alpha - \mathbb{E}_{H_0}(\alpha), \beta - \mathbb{E}_{H_0}(\beta)) \text{COV}_{H_0}^{-1} \left(\begin{pmatrix} \alpha \\ \beta \end{pmatrix} \right) \begin{pmatrix} \alpha - \mathbb{E}_{H_0}(\alpha) \\ \beta - \mathbb{E}_{H_0}(\beta) \end{pmatrix}.$$

The GPK can be decomposed into $\text{GPK} = Z_W^2 + Z_D^2$, where Z_W and Z_D are the standardized versions (with expectation and variance under H_0) of

$$\begin{aligned}W &= \frac{n_1}{N} \alpha + \frac{n_2}{N} \beta \\ D &= n_1(n_1 - 1) \alpha - n_2(n_2 - 1) \beta.\end{aligned}$$

Based on this observation, they further generalize W to

$$W_r = r \frac{n_1}{N} \alpha + \frac{n_2}{N} \beta$$

and Z_W to $Z_{W,r}$. Fast tests based on $Z_{W,r}$ are proposed as the asymptotic distribution of $Z_W = Z_{W,1}$ is complicated, but that of $Z_{W,r}, r \neq 1$, is a standard normal under mild assumptions. One fast test fGPK uses the Bonferroni adjusted test result of the tests based on $Z_D, Z_{W,1.2} =: ZW_1$ and $Z_{W,0.8} =: ZW_2$, the other fast test fGPK_M uses the Bonferroni adjusted test result of the tests based on $Z_{W,1.2}$ and $Z_{W,0.8}$. For GPK (as well as for fGPK and fGPK_M), a permutation test can be performed.

The new implementation `GPK()` based on the `kertests()` function from the `kerTests` package (Song and Chen 2023c) performs by default the fast test version instead of a permutation test, and the bandwidth parameter σ of the RBF kernel that is used as K is chosen via the median heuristic using the function `med_sigma()` of the `kerTests` package. The median heuristic sets the bandwidth of the kernel to the median value of all pairwise distances in the pooled sample (Sriperumbudur, Fukumizu, Gretton, Lanckriet, and Schölkopf 2009). When the fast test is performed, all three test statistics, ZW_1, ZW_2 , and Z_D are returned together with the asymptotic p value if `n.perm` = 0 or the permutation p value if `n.perm` > 0, respectively. For the GPK statistic, only the permutation test is available as its null distribution cannot be accessed. Therefore, if the

number of permutations is set to zero, the fast test is always performed. This holds even if `fast` is set to `FALSE` (with a warning).

5.14. LP test Mukhopadhyay and Wang (2020b)

For the test of Mukhopadhyay and Wang (2020b), a nonparametrically designed set of orthogonal functions (LP polynomials) is obtained by orthonormalizing a set of functions constructed as orthonormal polynomials of mid-distribution transforms. These are used for the construction of a polynomial kernel of degree 2 that encodes the similarity between two data points in the LP-transformed domain. The values of the kernel Gram matrix are then used as weights on a graph with the pooled sample as vertices. The idea is to cluster points for the graph into k groups that have higher connectivity and compare how closely related this clustering is to the true memberships of the k distributions. Then the problem reduces to testing independence, which can be accomplished by determining whether all of the LP comeans are zero.

The test is implemented in the `LPKsample` package (Mukhopadhyay and Wang 2020a). The new implementation offers the additional option to sum over all components instead of summing over the significant components only, which might be of interest when using the statistic as a data similarity measure without testing. By default, this is disabled (`sum.all = FALSE`). When only summing over the significant components, the returned test statistic is always equal to zero when no component is significant.

5.15. C2ST (Lopez-Paz and Oquab 2017)

The *classifier two-sample test* is already described in Section 1.3. For the C2ST, the classifier can be specified by the user and defaults to K -nearest neighbors. Possible options are all models accepted by `caret::train()`. For a list of classification models, call e.g.,

```
R> names(caret::getModelInfo())[sapply(caret::getModelInfo(), function(x) {
+   "Classification" %in% x$type
+ })]
```

5.16. Multisample cross-match tests of Mukherjee *et al.* (2022) and Petrie (2016)

The tests of Mukherjee *et al.* (2022) and Petrie (2016) generalize the Rosenbaum cross-match test to multiple samples by calculating the cross-counts for all pairs of samples based on the optimal non-bipartite matching on the pooled sample and taking the Mahalanobis distance or simply the sum of the cross-counts, respectively, as the test statistics. New functions `MMCM()` and `Petrie()` were implemented. There exist implementations of these methods in the R package `multicross` (Agarwal *et al.* 2020), but the package is archived on CRAN, and the implementation makes it impossible to access the test statistic and p value as numeric values. Therefore, here the functions were re-implemented from scratch. To ensure that the new functions are not derivations of the `multicross` versions, they were implemented by an author who had not looked at the `multicross` implementations before. The functions implement the formulas from Section 2 of Mukherjee *et al.* (2022). The new output is again of class `'htest'` and contains the test statistic value and the p value as a numeric value. The `nbpMatching` package (Beck, Lu, and Greevy 2024) is used for calculating the optimal non-bipartite matching. Note that in case of ties in the distance matrix, the optimal non-bipartite matching might not be defined uniquely. In the current implementation, the observations in the pooled sample are ordered as supplied by the user. When searching for a match, the `nbpMatching` implementation of the optimal non-bipartite matching algorithm starts at the end of the pooled sample. Therefore, with many ties (e.g., for categorical data), observations from the first dataset are often matched with ones from the last dataset, and so on. This might affect the validity of the test negatively since, even under the null, more cross counts than expected are observed. A random ordering of the pooled sample might help solve this issue, but would result in the observed test statistic value depending on this random ordering and is therefore not implemented.

5.17. Test using the shortest Hamiltonian path (Biswas *et al.* 2014)

Biswas *et al.* (2014) suggest a graph-based test similar to those of Friedman and Rafsky (1979) and Rosenbaum (2005) but using the shortest Hamiltonian path as the similarity graph. Since calculating the Hamiltonian path is an NP hard problem, the implementation of `BMG()` is based on Kruskal's algorithm, which is a heuristic

approach to find the shortest Hamilton Path within the pooled dataset as suggested in Biswas *et al.* (2014). Here, it is implemented as follows:

1. Create an edge list of the fully connected graph on the pooled sample, sorted by increasing Euclidean distance of the corresponding vertices.
2. For each edge, check if (i) an addition of this edge leads to a cyclic graph (using `IsAcyclic()` from the **rlemon** package (Agarwal, Tewari, and Errickson 2023)) and (ii) an addition of this edge leads to a degree larger than two in any (used) vertex. If both criteria are not met, keep the corresponding edge.
3. Return the reduced edge list, containing only edges needed to construct the Hamilton path.

For pooled sample sizes $N < 1030$, an exact test can be performed. For $N \geq 1030$ calculation of the exact runs statistic cannot be performed due to terms involved in the calculation becoming too large for representing them as floating point numbers in R. In the exact case, the p values using the null distribution of the univariate runs statistic (Biswas *et al.* 2014) are calculated. If an asymptotic test is performed, the asymptotic null distribution is used instead.

5.18. Rank Energy statistic (Deb and Sen 2021)

The test of Deb and Sen (2021) is a rank version of the Energy statistic. The multivariate ranks are assigned using optimal transport. The implementation is based on R code for the paper (<https://github.com/NabarunD/MultiDistFree>). It wraps up tidied-up versions of the `computestatistic()` and `gensamdist()` given there. The implementation uses the **randtoolbox** package (Christophe and Petr 2024) for random number generation, the **clue** package (Hornik 2005, 2024) to solve the assignment problem for ranking, and the **energy** package (Rizzo and Szekely 2024) for implementation of the Energy statistic.

5.19. Decision tree-based dataset similarity: Ganti *et al.* (1999) and Ntoutsis *et al.* (2008)

The methods of Ganti *et al.* (1999) and Ntoutsis *et al.* (2008) work by determining the partition induced by a decision tree fit to each dataset and then intersecting these partitions and calculating certain probability estimates on the resulting intersection. A description of the method of Ntoutsis *et al.* (2008) is given in Section 1.3. Ganti *et al.* (1999) calculate a decision tree model for each of the two datasets and calculate the greatest common refinement (GCR) induced by these trees. That is the intersection of the partitions of the sample space induced by each tree. A visualization of the computation of the GCR is given in Figure 2. Ganti *et al.* (1999) then compare the distribution of both datasets over this GCR. Let n_r denote the number of segments of the GCR, p_i the proportion of observations of $X^{(1)}$ that map to the i -th segment, and q_i the respective proportion of observations of $X^{(2)}$ mapping to the i -th segment. Then Ganti *et al.* (1999) compare the vector p and q by a difference function $f : \mathbb{R}^{n_r} \rightarrow \mathbb{R}^{n_r}$ and aggregate the results from that by an aggregate function $g : \mathbb{R}^{n_r} \rightarrow \mathbb{R}$ to obtain a measure of distance between the two datasets

$$\text{GAN} = g(f(p, q)).$$

Large values then indicate differences between the datasets. They propose the absolute difference function

$$f_a(p, q)_i = |p_i - q_i|,$$

and the scaled difference function

$$f_s(p, q)_i = \begin{cases} \frac{|p_i - q_i|}{(p_i + q_i)/2}, & \text{if } (p_i + q_i) > 0, i = 1, \dots, n_r. \\ 0, & \text{otherwise} \end{cases}.$$

For the aggregate function, they propose the sum or maximum of the values from the difference function. For using the sum as the aggregate function together with either f_a or f_s , it can be shown that the GCR is optimal in the sense that it gives the lowest value over all common refinements. For using the maximum, this property is not fulfilled. Ganti *et al.* (1999) propose using a Bootstrap test procedure for assessing whether or not the two datasets are generated by the same data-generating process.

We use the **rpart** package (Therneau and Atkinson 2025) for tree estimation. In the frame of a tree object fit with `rpart()`, the nodes are numbered starting with 1 at the root, following the rule that the left child node gets the ID of the parent times 2 and the right child node gets the ID of the parent times 2 plus 1. This allows

us to easily trace back the decision rules from a leaf node to the root using integer division by 2. Moreover, the split rules can be easily accessed using the `labels()` function on the tree object. We iterate over leaves and collect all split rules on each path from the leaf to the root. Suppose no upper or lower limit is specified by any split rule for a certain variable in this way. In that case, we set this limit to the minimum or maximum, respectively, of this variable over both datasets. This ensures that each observation in any of the two datasets falls into some part of the intersected partition later on. The resulting set of ranges for all variables for each leaf node gives us the partition induced by the tree. The resulting partitions are intersected as described in Ganti *et al.* (1999) and Ntoutsis *et al.* (2008). For Ntoutsis *et al.* (2008), all three methods presented in the original article (see also Section 1.3) are implemented. No test is performed. For Ganti *et al.* (1999), the difference and aggregation functions can be supplied by the users. The suggested choices f_a and f_s , i.e., taking the absolute differences between the joint probabilities calculated on the GCR or normalizing this difference with the sum of both probabilities, are readily implemented. The default different function is set to f_a , and the default aggregation function is set to the sum. A permutation test can be performed.

Neither Ntoutsis *et al.* (2008) nor Ganti *et al.* (1999) discuss the hyperparameter choice for the decision trees. Here, we offer the options to use the default parameter settings of `rpart()` or to tune the hyperparameters. For tuning the hyperparameters, we use the `best.rpart()` function of the **e1071** package (Meyer, Dimitriadou, Hornik, Weingessel, and Leisch 2024). The parameters `minsplit`, `minbucket`, and `cp` of the tree can be tuned. The ranges that are used here for tuning are chosen based on (Bischl, Binder, Lang, Pielok, Richter, Coors, Thomas, Ullmann, Becker, Boulesteix, Deng, and Lindauer 2021). Tuning is enabled by default but can be disabled by setting `tune` to `FALSE`. Cross-validation is used for tuning. The number of evaluations (`n.eval`) is set to 100 as a default, and the number of folds (`k`) is set to 5. Both values can be customized by the user. The remaining calculation works the same for a tuned or untuned tree model.

By default, the number of permutations is set to 0, corresponding to not performing any test. An implementation for categorical data for the method of Ganti *et al.* (1999) is also supplied. This comes with the following difficulties. If a category is only observed in one dataset and not in the other, or even if just not all combinations of categories are observed, it might happen that at a certain split, not all levels of the respective variable are observed in the remaining data at that split. Then it is unclear which child node the missing level gets assigned to. In the `rpart::rpart()` implementation that we use here, the label does not get assigned at all. If now in the other dataset, the combination with this label is present, the respective data points do not fit anywhere in the intersected partition. Therefore, the calculated probabilities in the joint distribution do not sum to one anymore. In these cases, a warning is printed. It might still give a useful measure of dataset distance, but the interpretation and theoretical results might not hold anymore. Also note that for deep trees, the intersection in practice often reduces to all combinations of categories of the variables. Therefore, the measure reduces to the differences in frequency of all category combinations in these cases, but is far more complicated and time-consuming to calculate.

5.20. OTDD (Alvarez-Melis and Fusi 2020)

A description of the optimal transport dataset distance can be found in Section 1.3. There is a Python implementation of the method (<https://github.com/microsoft/otdd>) that was used as a rough orientation here. Compared to that, the JDOT option is deprecated. The new implementation uses the Wasserstein distance implementation from **approxOT** package (Dunipace 2024) and the matrix square root from **expm** package (Maechler, Dutang, and Goulet 2024). Note that the solution of the optimal transport between two distributions is given by their q -Wasserstein distance to the power of q . There are different options for the method to calculate the optimal transport-based dataset distance. First case: chosen method is `"augmentation"`. In this case, the variable means and the covariance matrix of each dataset, reduced to each target observation value in that dataset, are calculated. The mean vector and the vectorized covariance matrix (column-wise) corresponding to the target value are appended to each observation in each dataset. Then, the q -Wasserstein distance to the power of q of these augmented datasets is calculated. Note that this calculation assumes commuting covariance matrices of all label distributions (rarely fulfilled in practice) and that the feature space metric coincides with the ground cost of the optimal transport problem on the labels (Alvarez-Melis and Fusi 2020). Second case: chosen method is `"precomputed.labeldist"`. In this case, both the distance matrix for the label distributions and the distance matrix for the features are calculated, and the corresponding distances are added with weights `lambda.x` and `lambda.y`, respectively, to calculate a cost matrix of all observations. In case of `sinkhorn = FALSE`, i.e., for the exact calculation, only the costs from each observation from the first dataset to each observation from the second dataset are needed. In the case of using debiased Sinkhorn approximation, additionally, the costs within each dataset are needed. For calculating

the distance matrices of the label distributions, there are again different options:

1. `inner.ot.method = "exact"`. The Wasserstein distance for each label pair between the datasets is reduced to the observations where the target value equals the corresponding label is calculated. There are options for using the (debiased) Sinkhorn approximation and changing the parameters of the Wasserstein distance and the ground cost metric.
2. `inner.ot.method = "gaussian.approx"`. The label distributions are approximated by Gaussians, which leads to a simple closed-form solution of the optimal transport problem that uses only the means and covariances. The calculation includes calculating multiple matrix square roots of covariance matrices, which might get costly if the number of variables is high. Moreover, this calculation fails if the estimated covariance matrix is not numerically positive definite. This might happen especially for $N < p$ settings.
3. `inner.ot.method = "only.means"`. The former is further simplified by using only the means (i.e., assuming equal covariance matrices in all label distributions).
4. `inner.ot.method = "naive.upperbound"`. A distribution-agnostic upper bound for the optimal transport between the label distributions is calculated that again only relies on the means and covariance matrices of these distributions.

5.21. Jeffreys Divergence

Jeffreys divergence (Jeffreys 1997) is the symmetrized version

$$J(F_1, F_2) = \text{KL}(F_1, F_2) + \text{KL}(F_2, F_1),$$

of the Kullback Leibler (KL) divergence (Kullback and Leibler 1951)

$$\text{KL}(F_1, F_2) := \int \log \left(\frac{f_1(x)}{f_2(x)} \right) f_1(x) dx.$$

Within the `Jeffreys()` function, Jeffreys divergence is calculated as the sum of the two KL-divergences (Kullback and Leibler 1951) where each dataset is used as the first once. The KL-divergences are calculated using density ratio estimation as recommended in Sugiyama, Liu, du Plessis, Yamanaka, Yamada, Suzuki, and Kanamori (2013). For this, the `densratio()` function from the `densratio` package (Makiyama 2019) is used. By default, the method `KLIEP` is chosen as suggested by Sugiyama *et al.* (2013). The `densratio` package was preferred here over the alternative package `densityratio` (Volker 2024) as it is available on CRAN.

5.22. Biswas and Ghosh (2014)

The statistic of Biswas and Ghosh (2014) uses inter-point distances and is defined as

$$T = \|\hat{\mu}_{D_{F_1}} - \hat{\mu}_{D_{F_2}}\|_2^2, \text{ where}$$

$$\hat{\mu}_{D_{F_1}} = \left[\hat{\mu}_{F_1 F_1} = \frac{2}{n_1(n_1 - 1)} \sum_{i=1}^{n_1} \sum_{j=i+1}^{n_1} \|X_i^{(1)} - X_j^{(1)}\|, \hat{\mu}_{F_1 F_2} = \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \|X_i^{(1)} - X_j^{(2)}\| \right]^\top,$$

$$\hat{\mu}_{D_{F_2}} = \left[\hat{\mu}_{F_1 F_2} = \frac{1}{n_1 n_2} \sum_{i=1}^{n_1} \sum_{j=1}^{n_2} \|X_i^{(1)} - X_j^{(2)}\|, \hat{\mu}_{F_2 F_2} = \frac{2}{n_2(n_2 - 1)} \sum_{i=1}^{n_2} \sum_{j=i+1}^{n_2} \|X_i^{(2)} - X_j^{(2)}\| \right]^\top.$$

For testing, the scaled statistic

$$T^* = \frac{N\hat{\lambda}(1 - \hat{\lambda})}{2\hat{\sigma}_0^2} T \text{ with}$$

$$\hat{\lambda} = \frac{n_1}{N},$$

$$\hat{\sigma}_0^2 = \frac{n_1 S_1 + n_2 S_2}{N}, \text{ where}$$

$$S_1 = \frac{1}{\binom{n_1}{3}} \sum_{1 \leq i < j < l \leq n_1} \|X_i^{(1)} - X_j^{(1)}\| \cdot \|X_i^{(1)} - X_l^{(1)}\| - \hat{\mu}_{F_1 F_1}^2 \text{ and}$$

$$S_2 = \frac{1}{\binom{n_2}{3}} \sum_{1 \leq i < j < l \leq n_2} \|X_i^{(2)} - X_j^{(2)}\| \cdot \|X_i^{(2)} - X_l^{(2)}\| - \hat{\mu}_{F_2 F_2}^2$$

is used as it is asymptotically χ_1^2 -distributed. The new function `BG2014()` implements the [Biswas and Ghosh \(2014\)](#) test from scratch. `stats::dist()` is used to calculate the Euclidean distance matrix on the pooled sample. The statistic T and the scaled test statistic T^* are implemented according to the formulas above. A permutation test is implemented by permuting the distance matrix, recalculating the test statistic T for the permuted distances, and calculating the p value as the proportion of permuted test statistics larger than the observed test statistic. An asymptotic test is implemented using the asymptotic result from Theorem 4.1 of [Biswas and Ghosh \(2014\)](#), i.e., calculating the p value as `stats::pchisq(T*, lower.tail = FALSE)`.

5.23. Engineer metric

The L_q -Engineer metric is defined as

$$\text{EN}(X, Y; q) = \left[\sum_{i=1}^p |\mathbb{E}(X_i) - \mathbb{E}(Y_i)|^q \right]^{\min(q, 1/q)} \quad \text{with } q > 0,$$

where X_i, Y_i denote the i th component of the p -dimensional random vectors $X \sim F_1$ and $Y \sim F_2$. A new function `engineerMetric()` is implemented. Since the Engineer metric is simply the L_q -distance of the expectations of two random vectors, it is estimated as the L_q -distance of the column means of the datasets. For the distance calculation, the base function `norm()` is used, and different options for the L_q norm are available via the `type` argument.

5.24. Schilling (1986) and Henze (1988) test

The Schilling-Henze test uses the mean within-sample edge-count, i.e.,

$$\text{SH} := L := \frac{1}{KN}(R_1 + R_2)$$

in a K -nearest neighbor graph as the test statistic. It is implemented from scratch as follows.

1. Calculate K -nearest neighbor (NN) edge matrix on the pooled sample (distance function returning a distance matrix and K are inputs of the function), i.e. create a matrix where the first column is each observation number repeated K times, and the second column are the corresponding K nearest neighbors of that observation. For the calculation of the K -NN graph, a function can be supplied by the user. Pre-implemented options include a wrapper for the `knn()` function from the `dbscan` package ([Hahsler, Piekenbrock, and Doran 2019](#)) and the fast (approximative) K -NN algorithm implemented in the `get.knn()` function from the `FNN` package ([Beygelzimer, Kakadet, Langford, Arya, Mount, and Li 2024](#)).
2. Count the number of rows where both observations come from the same sample L (i.e., either both have observation number $\leq n_1$ or both have observation number $> n_1$)
3. Calculate the quantities $\mathbb{E}_{H_0}(L)$ and $\text{VAR}_{H_0}(L)$ from proposition 2.1 in [Henze \(1988\)](#)
4. Calculate the standardized test statistic $L^* = (L - \mathbb{E}_{H_0}(L)) / \sqrt{\text{VAR}_{H_0}(L)}$
5. When performing a permutation test, permute the distance matrix on the pooled sample, recalculate L , and calculate the proportion of permuted test statistics that are larger than the observed value of L
6. When performing an asymptotic test, use the asymptotic normal distribution of Z as proposed in Remark 5.1 of [Henze \(1988\)](#).
7. The observed value of L^* is returned in the result as the `statistic`, the observed L is returned as the `estimate`.

The default for K is set to one. This is rather arbitrary based on computational speed, as there is no good rule for choosing K so far proposed in the literature ([Aslan and Zech 2005](#)).

5.25. Barakat et al. (1996) Generalization of the Schilling-Henze Test

[Barakat et al. \(1996\)](#) generalize the Schilling-Henze nearest neighbor test to circumvent choosing the number of nearest neighbors. Their test statistic is the sum of edge counts for all values of K for the K -nearest neighbor graph. The resulting test is equivalent to a sum of Wilcoxon rank sums. It requires samples in the Euclidean

space \mathbb{R}^p and it is assumed that there are no ties in ranking w.r.t. to nearness.

Within our implementation, we do not explicitly calculate the K -nearest neighbor graph for all possible values of K as this would be highly inefficient. Instead, the distance matrix on the pooled sample is calculated with a user-specified distance function (Euclidean distance calculated via `stats::dist()` by default), and the column-wise orderings of the distances, excluding the diagonal elements, are calculated. Then, the cumulative numbers of the elements smaller than n_1 are calculated for the first n_1 columns of the orderings, corresponding to the numbers of within-sample edges in the first sample in the K -nearest neighbor graph for $K = 1, \dots, N - 1$. Analogously, the cumulative numbers of the elements greater than n_1 are calculated for the remaining n_2 columns of the orderings, corresponding to the numbers of within-sample edges in the second sample in the K -nearest neighbor graph for $K = 1, \dots, N - 1$. Lastly, all these cumulative numbers are summed up, which corresponds to the Barakat *et al.* (1996) test statistic. A permutation test is implemented using the `boot::boot()` function. For that, the distances are permuted directly, and the calculation is repeated for the permuted distance matrix, which circumvents the costly recalculation of the distances for each permutation.

5.26. Tree-based test (Yu *et al.* 2007)

Yu *et al.* (2007) propose a permutation test that uses the classification error of a classification tree that distinguishes between the two datasets. The implementation of the test is based on the `C2ST()` function as the methods work very similarly. Here, we set the classifier to "rpart", i.e., a CART. Instead of the classification accuracy as for the C2ST, the classification error, i.e., $1 - \text{Accuracy}$, is returned. A permutation test is implemented using the `boot::boot()` framework, and the permutation p value is calculated as the proportion of the number + 1 of permuted test statistics smaller than the observed value divided by the number of permutations. Yu *et al.* (2007) do not propose any asymptotic test, but since their test fits into the framework of Lopez-Paz and Oquab (2017), the binomial test proposed there and implemented in the `Ecume::classifier_test()` function utilized by `C2ST()` is still valid and therefore kept in the implementation.

5.27. Characteristic distance (Li *et al.* 2022)

The characteristic distance is defined as

$$\begin{aligned} \text{CD}(X, Y) = & \mathbb{E} [\| \mathbb{E} (\exp (i \langle X'', X - X' \rangle) | X - X') \\ & - \mathbb{E} (\exp (i \langle Y, X - X' \rangle) | X - X') \|^2] \\ & + \mathbb{E} [\| \mathbb{E} (\exp (i \langle X, Y - Y' \rangle) | Y - Y') \\ & - \mathbb{E} (\exp (i \langle Y'', Y - Y' \rangle) | Y - Y') \|^2], \end{aligned}$$

where X', X'' and Y', Y'' denote independent copies of $X \sim F_1$ and $Y \sim F_2$, respectively. An empirical version is obtained by replacing the conditional expectations with empirical means. The implementation calculates the empirical characteristic distance between two datasets. For both summands, Euler's formula is used for every entry of the inner product defined in Li *et al.* (2022). Both mean values are calculated, and the squared complex modulus of the difference between both means is calculated. Since the inner product leads to a symmetric matrix, only an upper triangular matrix is calculated, and the final sum is multiplied by two. A permutation test with `n.perm` permutations and random seed `seed` for reproducibility is performed.

5.28. Constrained Minimum Distance (Tatti 2007)

The *constrained minimum (CM) distance* uses a *feature function* $S : \mathcal{X} \rightarrow \mathbb{R}^m$ that maps points from the sample space \mathcal{X} to a real vector. The *frequency* $\theta \in \mathbb{R}^m$ of S with respect to dataset $X^{(j)}$ is the average of the values of S

$$\theta_j = \frac{1}{N} \sum_{i=1}^{n_i} S(X_i^{(j)}), j = 1, 2.$$

The CM distance is then defined as

$$D_{\text{CM}}(X^{(1)}, X^{(2)} | S)^2 = (\theta_1 - \theta_2)^\top \text{COV}^{-1}(S) (\theta_1 - \theta_2)$$

with

$$\text{COV}(S) = \frac{1}{|\mathcal{X}|} \sum_{\omega \in \mathcal{X}} S(\omega) S(\omega)^\top - \left(\frac{1}{|\mathcal{X}|} \sum_{\omega \in \mathcal{X}} S(\omega) \right) \left(\frac{1}{|\mathcal{X}|} \sum_{\omega \in \mathcal{X}} S(\omega) \right)^\top.$$

It has to be assumed that the feature space \mathcal{X} is finite and can be enumerated. For binary data and S chosen as the conjunction function, i.e., S is one if all components of an observation are one, and zero otherwise, or as the parity function, i.e. S is one if an odd number of components of an observation are one, and zero otherwise, the CM distance reduces to

$$D_{\text{CM}}(X^{(1)}, X^{(2)}|S) = 2\|\theta_1 - \theta_2\|_2.$$

This special case for binary data is implemented first. It includes the option to use either the means as features (example 3 in Tatti (2007)) or the means and covariances (example 4 in Tatti (2007)). Note that there is an error in the calculation of the covariance matrix in A.4 Proof of Lemma 8 in Tatti (2007). The correct covariance matrix has the form $\text{COV}[T_{\mathcal{F}}] = 0.25I$ since $\text{VAR}[T_A] = \text{E}[T_A^2] - \text{E}[T_A]^2 = 0.5 - 0.5^2 = 0.25$ following from the correct statement that $\text{E}[T_A^2] = \text{E}[T_A] = 0.5$. Therefore, formula (4) changes to $d_{\text{CM}}(D_1, D_2|S_{\mathcal{F}}) = 2\|\theta_1 - \theta_2\|_2$ and the formula in example 3 changes to $d_{\text{CM}}(D_1, D_2|S_1) = 2\|\theta_1 - \theta_2\|_2$. Our implementation is based on these corrected formulas. If the original formula was used, the results on the same data calculated with the formula for the binary special case and the results calculated with the general formula differ by a factor of $\sqrt{2}$. For the general case for categorical data, the user has to specify a feature function S mapping a point in the sample space to a real vector. Additionally, either the covariance matrix $\text{COV}[S]$ if known or the sample space has to be given. If both are given, the supplied covariance matrix is used and not recalculated. The constrained minimum distance is calculated using Theorem 1 in Tatti (2007), i.e., the formulas given above. Therefore, the supplied or calculated $\text{COV}[S]$, respectively, has to be invertible.

5.29. Biau and Györfi (2005)

Biau and Györfi (2005) test for homogeneity of two (multivariate) datasets by calculating the L_1 -distance between the two empirical distributions restricted to a finite partition. For this, a finite partition of the subspace spanned by the two datasets has to be defined. By default, we define a rectangular partition under the assumption of approximately equal cell probabilities. The number of elements of the partition m_n is chosen according to the convergence criteria in Biau and Györfi (2005) as $n^{0.8}$, where the exponent can be varied as an argument (`exponent`). For each dimension, $m_n^{1/p} + 1$ equidistant cut-points are created along the range of both datasets to define the partition. It must be ensured that there are at least three cut-points per dimension (min, max, and one point splitting the data into two bins). The argument `eps` ensures that the partition covers all data points by adding some small value to the data range. Alternative partition functions can be provided via the `partition` argument. After calculating the partition, all data points are assigned to an element of the partition along the defined cut-points. Last, the L_1 distance between the empirical distribution functions restricted to the elements of the partition is calculated.

5.30. DiProPerm test (Wei et al. 2016)

Wei et al. (2016) propose their *direction-projection-permutation* (*DiProPerm*) test for which a univariate two-sample statistic is applied to the projection of the datasets onto the normal vector of a separating hyperplane. For this, a linear classification method like a support vector machine (SVM) or distance weighted discrimination (DWD) is used to calculate such a separating hyperplane. A permutation test is then performed for the univariate statistic applied to the projection onto the normal vector. Possible options for the univariate statistic would be the mean difference, the two-sample t -statistic, or the area under the curve (AUC). There is an implementation in the **diproperm** package (Allmon et al. 2021) which is currently archived. Our implementation is independent of that implementation. It has the following advantages.

- All suggested univariate two-sample statistics from the paper, i.e., mean difference, t test statistic and AUC are implemented. Additional two-sample statistics can be used if a suitable function is supplied via the `stat.fun` argument.
- Additional binary linear classifiers other than the DWD and SVM suggested in the original paper can easily be used by supplying a suitable function via the `dipro.fun` argument.
- The results of the new function are reproducible by setting a random seed.
- The new implementation does not rely on global variables
- The p value is returned as numeric instead of character.
- The output is an object of class `'htest'` for pretty displaying of the results.

One restriction of the new function is that it no longer supports balanced permutation. That was necessary to ensure the reproducibility, which we considered a trade-off worth making since the use of balanced permutation is controversial anyway, see [Southworth, Kim, and Owen \(2009\)](#), and reproducibility is essential for permutation tests.

Acknowledgments

This work has been supported (in part) by the Research Training Group “Biostatistical Methods for High-Dimensional Data in Toxicology” (RTG 2624, Project P1) funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation - Project Number 427806116).

We would like to thank Nabarun Deb and Bodhisattva Sen for allowing us to use their R implementation of their test for the package. Moreover, we would like to thank David Alvarez-Melis, whose Python implementation of the OTDD was the basis for our R implementation.

References

- Agarwal A, Tewari A, Errickson J (2023). *rlemon: R Access to LEMON Graph Algorithms*. R package version 0.2.1, URL <https://CRAN.R-project.org/package=rlemon>.
- Agarwal SMD, Bhattacharya B, Zhang NR (2020). *multicross: A Graph-Based Test for Comparing Multivariate Distributions in the Multi Sample Framework*. R package version 2.1.0, URL <https://CRAN.R-project.org/package=multicross>.
- Allmon AG, Marron J, Hudgens MG (2021). *diproperm: Conduct Direction-Projection-Permutation Tests and Display Plots*. R package version 0.2.0, URL <https://CRAN.R-project.org/package=diproperm>.
- Alvarez-Melis D, Fusi N (2020). “Geometric Dataset Distances via Optimal Transport.” In *Advances in Neural Information Processing Systems*, volume 33, pp. 21428–21439. Curran Associates, Inc.
- Alvarez-Melis D, Fusi N (2020). “Measuring dataset similarity using optimal transport.” URL <https://www.microsoft.com/en-us/research/blog/measuring-dataset-similarity-using-optimal-transport/>.
- Aslan B, Zech G (2005). “New Test for the Multivariate Two-Sample Problem Based on the Concept of Minimum Energy.” *Journal of Statistical Computation and Simulation*, **75**(2), 109–119. ISSN 0094-9655. doi:10.1080/00949650410001661440.
- Bahr R (1996). *Ein neuer Test für das mehrdimensionale Zwei-Stichproben-Problem bei allgemeiner Alternative*. Ph.D. thesis, Universität Hannover.
- Barakat AS, Quade D, Salama IA (1996). “Multivariate Homogeneity Testing Using an Extended Concept of Nearest Neighbors.” *Biometrical Journal*, **38**(5), 605–612. ISSN 1521-4036. doi:10.1002/bimj.4710380509.
- Baringhaus L, Franz C (2004). “On a New Multivariate Two-Sample Test.” *Journal of Multivariate Analysis*, **88**(1), 190–206. ISSN 0047-259X. doi:10.1016/S0047-259X(03)00079-4.
- Baringhaus L, Franz C (2010). “Rigid Motion Invariant Two-Sample Tests.” *Statistica Sinica*, **20**(4), 1333–1361. ISSN 1017-0405.
- Beck C, Lu B, Greevy R (2024). *nbpMatching: Functions for Optimal Non-Bipartite Matching*. R package version 1.5.6, URL <https://CRAN.R-project.org/package=nbpMatching>.
- Beygelzimer A, Kakadet S, Langford J, Arya S, Mount D, Li S (2024). *FNN: Fast Nearest Neighbor Search Algorithms and Applications*. R package version 1.1.4, URL <https://CRAN.R-project.org/package=FNN>.
- Biau G, Györfi L (2005). “On the Asymptotic Properties of a Nonparametric L_1 -Test Statistic of Homogeneity.” *IEEE Transactions on Information Theory*, **51**(11), 3965–3973. ISSN 1557-9654. doi:10.1109/TIT.2005.856979.

- Bischl B, Binder M, Lang M, Pielok T, Richter J, Coors S, Thomas J, Ullmann T, Becker M, Boulesteix AL, Deng D, Lindauer M (2021). “Hyperparameter Optimization: Foundations, Algorithms, Best Practices and Open Challenges.” *arXiv:2107.05847 [cs, stat]*.
- Biswas M, Ghosh AK (2014). “A Nonparametric Two-Sample Test Applicable to High Dimensional Data.” *Journal of Multivariate Analysis*, **123**, 160–171. ISSN 0047-259X. doi:10.1016/j.jmva.2013.09.004.
- Biswas M, Mukhopadhyay M, Ghosh AK (2014). “A Distribution-Free Two-Sample Run Test Applicable to High-Dimensional Data.” *Biometrika*, **101**(4), 913–926. ISSN 0006-3444. doi:10.1093/biomet/asu045.
- Chen H, Chen X, Su Y (2018). “A Weighted Edge-Count Two-Sample Test for Multivariate and Object Data.” *Journal of the American Statistical Association*, **113**(523), 1146–1155. ISSN 0162-1459. doi:10.1080/01621459.2017.1307757.
- Chen H, Friedman JH (2017). “A New Graph-Based Two-Sample Test for Multivariate and Object Data.” *Journal of the American Statistical Association*, **112**(517), 397–409. ISSN 0162-1459. doi:10.1080/01621459.2016.1147356.
- Chen H, Zhang J (2017). *gTests: Graph-Based Two-Sample Tests*. R package version 0.2, URL <https://CRAN.R-project.org/package=gTests>.
- Chen H, Zhang NR (2013). “Graph-Based Tests for Two-Sample Comparisons of Categorical Data.” *Statistica Sinica*, **23**(4), 1479–1503. ISSN 1017-0405.
- Chen L, Dou WW, Qiao Z (2013). “Ensemble Subsampling for Imbalanced Multivariate Two-Sample Tests.” *Journal of the American Statistical Association*, **108**(504), 1308–1323. ISSN 0162-1459. doi:10.1080/01621459.2013.800763.
- Christophe D, Petr S (2024). *randtoolbox: Generating and Testing Random Numbers*. R package version 2.0.5.
- Chu L, Chen H (2019). “Asymptotic Distribution-Free Change-Point Detection for Multivariate and Non-Euclidean Data.” *The Annals of Statistics*, **47**(1), 382–414. ISSN 0090-5364, 2168-8966. doi:10.1214/18-AOS1691.
- Deb N, Sen B (2021). “Multivariate Rank-Based Distribution-Free Nonparametric Testing Using Measure Transportation.” *Journal of the American Statistical Association*, **118**(541), 1–16. ISSN 0162-1459. doi:10.1080/01621459.2021.1923508.
- Dray S, Dufour AB (2007). “The ade4 Package: Implementing the Duality Diagram for Ecologists.” *Journal of Statistical Software*, **22**(4), 1–20. doi:10.18637/jss.v022.i04.
- Dunipace EA (2024). *approxOT: approximate optimal transport*. R package version 1.1, URL <https://github.com/ericdunipace/approxOT>.
- Fisher RA (1936). “The Use of Multiple Measurements in Taxonomic Problems.” *Annals of Eugenics*, **7**(2), 179–188. ISSN 2050-1439. doi:10.1111/j.1469-1809.1936.tb02137.x. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x>.
- Franz C (2024). *cramer: Multivariate Nonparametric Cramer-Test for the Two-Sample-Problem*. R package version 0.9-4, URL <https://CRAN.R-project.org/package=cramer>.
- Friedman JH, Rafsky LC (1979). “Multivariate Generalizations of the Wald-Wolfowitz and Smirnov Two-Sample Tests.” *The Annals of Statistics*, **7**(4), 697–717. ISSN 0090-5364.
- Fukumizu K, Bach FR, Jordan MI (2004). “Dimensionality Reduction for Supervised Learning with Reproducing Kernel Hilbert Spaces.” *Journal of Machine Learning Research*, **5**, 73–99.
- Ganti V, Gehrke J, Ramakrishnan R, Loh WY (1999). “A Framework for Measuring Changes in Data Characteristics.” In *Proceedings of the 18th Symposium on Principles of Database Systems*, pp. 126–137.
- Gretton A, Borgwardt K, Rasch M, Schölkopf B, Smola A (2006). “A Kernel Method for the Two-Sample Problem.” In *Advances in Neural Information Processing Systems*, volume 19. MIT Press.

- Hahsler M, Piekenbrock M, Doran D (2019). “dbscan: Fast Density-Based Clustering with R.” *Journal of Statistical Software*, **91**(1), 1–30. doi:10.18637/jss.v091.i01.
- Hediger S, Michel L, Naef J (2024). *hypoRF: Random Forest Two-Sample Tests*. R package version 1.0.1, URL <https://CRAN.R-project.org/package=hypoRF>.
- Hediger S, Michel L, Näf J (2022). “On the Use of Random Forest for Two-Sample Testing.” *Computational Statistics & Data Analysis*, **170**, 107435. ISSN 0167-9473. doi:10.1016/j.csda.2022.107435. URL <https://www.sciencedirect.com/science/article/pii/S0167947322000159>.
- Heller R, Small D, Rosenbaum P (2024). *crossmatch: The Cross-Match Test*. R package version 1.4-0, URL <https://CRAN.R-project.org/package=crossmatch>.
- Henze N (1988). “A Multivariate Two-Sample Test Based on the Number of Nearest Neighbor Type Coincidences.” *The Annals of Statistics*, **16**(2), 772–783. ISSN 0090-5364.
- Hill AA, LaPan P, Li Y, Haney S (2007). “Impact of Image Segmentation on High-Content Screening Data Quality for SK-BR-3 Cells.” *BMC Bioinformatics*, **8**(1), 340. ISSN 1471-2105. doi:10.1186/1471-2105-8-340. URL <https://doi.org/10.1186/1471-2105-8-340>.
- Hornik K (2005). “A CLUE for CLUster Ensembles.” *Journal of Statistical Software*, **14**(12). doi:10.18637/jss.v014.i12.
- Hornik K (2024). *clue: Cluster Ensembles*. R package version 0.3-66, URL <https://CRAN.R-project.org/package=clue>.
- Huang Z (2022). *KMD: Kernel Measure of Multi-Sample Dissimilarity*. R package version 0.1.0, URL <https://CRAN.R-project.org/package=KMD>.
- Huang Z, Sen B (2023). “A Kernel Measure of Dissimilarity between M Distributions.” *Journal of the American Statistical Association*, pp. 1–27. ISSN 0162-1459. doi:10.1080/01621459.2023.2298036. URL <https://doi.org/10.1080/01621459.2023.2298036>.
- Jeffreys H (1997). “An Invariant Form for the Prior Probability in Estimation Problems.” *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences*, **186**(1007), 453–461. doi:10.1098/rspa.1946.0056.
- Karatzoglou A, Smola A, Hornik K, Zeileis A (2004). “kernlab – An S4 Package for Kernel Methods in R.” *Journal of Statistical Software*, **11**(9), 1–20. doi:10.18637/jss.v011.i09.
- Kuhn, Max (2008). “Building Predictive Models in R Using the caret Package.” *Journal of Statistical Software*, **28**(5), 1–26. doi:10.18637/jss.v028.i05. URL <https://www.jstatsoft.org/index.php/jss/article/view/v028i05>.
- Kullback S, Leibler RA (1951). “On Information and Sufficiency.” *The Annals of Mathematical Statistics*, **22**(1), 79–86. ISSN 0003-4851, 2168-8990. doi:10.1214/aoms/1177729694.
- Li X, Hu W, Zhang B (2022). “Measuring and Testing Homogeneity of Distributions by Characteristic Distance.” *Statistical Papers*. ISSN 1613-9798. doi:10.1007/s00362-022-01327-7.
- Lopez-Paz D, Oquab M (2017). “Revisiting Classifier Two-Sample Tests.” In *International Conference on Learning Representations*. URL <https://openreview.net/forum?id=SJkXfE5xx>.
- Maechler M, Dutang C, Goulet V (2024). *expm: Matrix Exponential, Log, 'etc'*. R package version 1.0-0, URL <https://CRAN.R-project.org/package=expm>.
- Makiyama K (2019). *densratio: Density Ratio Estimation*. R package version 0.2.1, URL <https://CRAN.R-project.org/package=densratio>.
- Meyer D, Dimitriadou E, Hornik K, Weingessel A, Leisch F (2024). *e1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071), TU Wien*. R package version 1.7-16, URL <https://CRAN.R-project.org/package=e1071>.

- Muandet K, Fukumizu K, Sriperumbudur B, Schölkopf B (2017). “Kernel Mean Embedding of Distributions: A Review and Beyond.” *Foundations and Trends® in Machine Learning*, **10**(1-2), 1–141. ISSN 1935-8237, 1935-8245. doi:10.1561/22000000060.
- Mukherjee S, Agarwal D, Zhang NR, Bhattacharya BB (2022). “Distribution-Free Multisample Tests Based on Optimal Matchings With Applications to Single Cell Genomics.” *Journal of the American Statistical Association*, **117**(538), 627–638. ISSN 0162-1459. doi:10.1080/01621459.2020.1791131.
- Mukhopadhyay S, Wang K (2020a). *LPKsample: LP Nonparametric High Dimensional K-Sample Comparison*. R package version 2.1, URL <https://CRAN.R-project.org/package=LPKsample>.
- Mukhopadhyay S, Wang K (2020b). “A Nonparametric Approach to High-Dimensional k-Sample Comparison Problems.” *Biometrika*, **107**(3), 555–572. ISSN 0006-3444. doi:10.1093/biomet/asaa015.
- Ntoutsi I, Kalousis A, Theodoridis Y (2008). “A General Framework for Estimating Similarity of Datasets and Decision Trees: Exploring Semantic Similarity of Decision Trees.” In *Proceedings of the 2008 SIAM International Conference on Data Mining (SDM)*, pp. 810–821. Society for Industrial and Applied Mathematics. ISBN 978-0-89871-654-2. doi:10.1137/1.9781611972788.73.
- Pan W, Tian Y, Wang X, Zhang H (2018). “Ball Divergence: Nonparametric Two Sample Test.” *The Annals of Statistics*, **46**(3), 1109–1137. ISSN 0090-5364. doi:10.1214/17-AOS1579.
- Paul B, De SK, Ghosh AK (2022a). *HDLSSkST: Distribution-Free Exact High Dimensional Low Sample Size k-Sample Tests*. R package version 2.1.0, URL <https://CRAN.R-project.org/package=HDLSSkST>.
- Paul B, De SK, Ghosh AK (2022b). “Some Clustering-Based Exact Distribution-Free k-Sample Tests Applicable to High Dimension, Low Sample Size Data.” *Journal of Multivariate Analysis*, **190**, 104897. ISSN 0047-259X. doi:10.1016/j.jmva.2021.104897. URL <https://www.sciencedirect.com/science/article/pii/S0047259X21001743>.
- Petrie A (2016). “Graph-Theoretic Multisample Tests of Equality in Distribution for High Dimensional Data.” *Computational Statistics & Data Analysis*, **96**, 145–158. ISSN 0167-9473. doi:10.1016/j.csda.2015.11.003.
- R Core Team (2024). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- Rizzo M, Székely G (2024). *energy: E-Statistics: Multivariate Inference via the Energy of Data*. R package version 1.7-12, URL <https://CRAN.R-project.org/package=energy>.
- Rizzo ML, Székely GJ (2010). “DISCO Analysis: A Nonparametric Extension of Analysis of Variance.” *The Annals of Applied Statistics*, **4**(2), 1034–1055. ISSN 1932-6157, 1941-7330. doi:10.1214/09-AOAS245.
- Rosenbaum PR (2005). “An Exact Distribution-Free Test Comparing Two Multivariate Distributions Based on Adjacency.” *Journal of the Royal Statistical Society B*, **67**(4), 515–530. ISSN 1369-7412.
- Roux de Bezieux H (2024). *Ecume: Equality of 2 (or k) Continuous Univariate and Multivariate Distributions*. R package version 0.9.2, URL <https://CRAN.R-project.org/package=Ecume>.
- Sarkar S, Ghosh AK (2020). “On Perfect Clustering of High Dimension, Low Sample Size Data.” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **42**(9), 2257–2272. ISSN 1939-3539. doi:10.1109/TPAMI.2019.2912599. Conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence, URL <https://ieeexplore.ieee.org/document/8695805>.
- Schilling MF (1986). “Multivariate Two-Sample Tests Based on Nearest Neighbors.” *Journal of the American Statistical Association*, **81**(395), 799–806. ISSN 0162-1459. doi:10.2307/2289012.
- Song H, Chen H (2022). “New Graph-Based Multi-Sample Tests for High-Dimensional and Non-Euclidean Data.” doi:10.48550/arXiv.2205.13787. ArXiv:2205.13787 [stat].

- Song H, Chen H (2023a). “Generalized Kernel Two-Sample Tests.” *Biometrika*, pp. 755–770. ISSN 1464-3510. doi:10.1093/biomet/asad068.
- Song H, Chen H (2023b). *gTestsMulti: New Graph-Based Multi-Sample Tests*. R package version 0.1.1, URL <https://CRAN.R-project.org/package=gTestsMulti>.
- Song H, Chen H (2023c). *kerTests: Generalized Kernel Two-Sample Tests*. R package version 0.1.4, URL <https://CRAN.R-project.org/package=kerTests>.
- Southworth LK, Kim SK, Owen AB (2009). “Properties of Balanced Permutations.” *Journal of Computational Biology: A Journal of Computational Molecular Cell Biology*, **16**(4), 625–638. ISSN 1557-8666. doi:10.1089/cmb.2008.0144.
- Sriperumbudur B, Fukumizu K, Gretton A, Lanckriet G, Schölkopf B (2009). “Kernel Choice and Classifiability for RKHS Embeddings of Probability Distributions.” In *Advances in Neural Information Processing Systems 22*, pp. 1750–1758. Max-Planck-Gesellschaft, Curran, Red Hook, NY, USA.
- Sriperumbudur BK, Gretton A, Fukumizu K, Lanckriet G, Schölkopf B (2008). “Injective Hilbert space embeddings of probability measures.” In *21st Annual Conference on Learning Theory (COLT 2008)*, pp. 111–122. Omnipress.
- Sriperumbudur BK, Gretton A, Fukumizu K, Schölkopf B, Lanckriet GRG (2010). “Hilbert Space Embeddings and Metrics on Probability Measures.” *Journal of Machine Learning Research*, **11**(50), 1517–1561. ISSN 1533-7928.
- Stolte M, Kappenberg F, Rahnenführer J, Bommert A (2024). “Methods for Quantifying Dataset Similarity: A Review, Taxonomy and Comparison.” *Statistics Surveys*, **18**, 163–298. ISSN 1935-7516. doi:10.1214/24-SS149.
- Sugiyama M, Liu S, du Plessis MC, Yamanaka M, Yamada M, Suzuki T, Kanamori T (2013). “Direct Divergence Approximation between Probability Distributions and Its Applications in Machine Learning.” *Journal of Computing Science and Engineering*, **7**(2), 99–111. ISSN 1976-4677. doi:10.5626/JCSE.2013.7.2.99.
- Sutherland JJ, Weaver DF (2004). “Three-Dimensional Quantitative Structure-Activity and Structure-Selectivity Relationships of Dihydrofolate Reductase Inhibitors.” *Journal of Computer-Aided Molecular Design*, **18**(5), 309–331. ISSN 0920-654X. doi:10.1023/b:jcam.0000047814.85293.da.
- Szabo A, Boucher K, Carroll WL, Klebanov LB, Tsodikov AD, Yakovlev AY (2002). “Variable selection and pattern recognition with gene expression data generated by the microarray technology.” *Mathematical Biosciences*, **176**(1), 71–98. ISSN 0025-5564. doi:10.1016/s0025-5564(01)00103-1.
- Székely GJ, Rizzo ML (2017). “The Energy of Data.” *Annual Review of Statistics and Its Application*, **4**(1), 447–479. doi:10.1146/annurev-statistics-060116-054026.
- Tatti N (2007). “Distances between Data Sets Based on Summary Statistics.” *Journal of Machine Learning Research*, **8**(1).
- Therneau T, Atkinson B (2025). *rpart: Recursive Partitioning and Regression Trees*. R package version 4.1.24, URL <https://CRAN.R-project.org/package=rpart>.
- Vaserstein LN (1969). “Markov Processes Over Denumerable Products of Spaces, Describing Large Systems of Automata.” *Problemy Peredachi Informatsii*, **5**(3), 64–72.
- Volker TB (2024). “densityratio: Distribution Comparison through Density Ratio Estimation.” doi:10.5281/zenodo.13881689. URL <https://github.com/thomvolker/densityratio>.
- Wei S, Lee C, Wichers L, Marron JS (2016). “Direction-Projection-Permutation for High-Dimensional Hypothesis Tests.” *Journal of Computational and Graphical Statistics*, **25**(2), 549–569. ISSN 1061-8600. doi:10.1080/10618600.2015.1027773.

- Yu K, Martin R, Rothman N, Zheng T, Lan Q (2007). “Two-sample Comparison Based on Prediction Error, with Applications to Candidate Gene Association Studies.” *Annals of Human Genetics*, **71**(1), 107–118. ISSN 1469-1809. doi:[10.1111/j.1469-1809.2006.00306.x](https://doi.org/10.1111/j.1469-1809.2006.00306.x).
- Zhang J, Chen H (2019). “Graph-Based Two-Sample Tests for Data with Repeated Observations.” doi:[10.48550/arXiv.1711.04349](https://doi.org/10.48550/arXiv.1711.04349). ArXiv:1711.04349 [stat].
- Zhu J, Pan W, Zheng W, Wang X (2021). “Ball: An R Package for Detecting Distribution Difference and Association in Metric Spaces.” *Journal of Statistical Software*, **97**(6), 1–31. doi:[10.18637/jss.v097.i06](https://doi.org/10.18637/jss.v097.i06).

Affiliation:

Marieke Stolte
 Department of Statistics
 TU Dortmund University
 Vogelpothsweg 87
 44227 Dortmund, Germany
 E-mail: stolte@statistik.tu-dortmund.de