

ITE in Python

(Information Theoretical Estimators - Release 1.1)

Zoltán Szabó

February 20, 2018

Contents

1	Introduction	1
2	Getting Started: Installation, Built-in Demos, Examples	3
3	Estimated Quantities and Estimators	5
3.1	Unconditional Quantities	5
3.1.1	Entropy	6
3.1.2	Mutual Information	6
3.1.3	Divergence	8
3.1.4	Association Measure	10
3.1.5	Cross Quantity	12
3.1.6	Kernel on Distributions	12
3.2	Conditional Quantities	13
3.2.1	Entropy	13
3.2.2	Mutual Information	13
3.3	Estimators	13
A	For Developers	18
A.1	Directory Structure	18
A.2	Running Doctests	18
A.3	Adding New Estimators	19
A.4	Parameter Passing for (Certain) Meta Estimators	19
B	Python ITE ↔ Matlab ITE	21
C	For Mathy People: Axioms of Concordance and Dependence	21

List of Examples

1	Entropy estimation	4
2	Mutual information estimation (association measure: similarly)	4
3	Divergence estimation (cross quantity, kernel on distributions: analogously)	5
4	Conditional entropy estimation	5
5	Conditional mutual information estimation	5

1 Introduction

Measuring the uncertainty, independence, association, inner product or distance of random variables is a central problem with numerous applications. Despite the large number of successful applications and emerging potentials, there are quite few available packages in the area which would help the comparison of different information theoretical measures and

estimation techniques.¹ To remedy this serious bottleneck and provide a platform in a rapidly evolving, free software environment we created the *Information Theoretical Estimators (ITE) in Python* toolbox. It

1. is the redesigned, Python implementation of the Matlab/Octave ITE toolbox².
2. can estimate numerous entropy, mutual information, divergence, association measures, cross quantities, and kernels on distributions (see the list below).
3. can be used to solve information theoretical optimization problems in a high-level way.
4. comes with several demos.

Some details:

- **Estimated quantities:**

- **entropy:** Shannon entropy, Rényi entropy, Tsallis entropy (Havrda and Charvát entropy), Sharma-Mittal entropy, Φ -entropy (f -entropy).
- **mutual information:** Shannon mutual information (total correlation, multi-information), Rényi mutual information, Tsallis mutual information, χ^2 mutual information (squared-loss mutual information, mean square contingency), L_2 mutual information, copula-based kernel dependency, kernel canonical correlation analysis (KCCA), kernel generalized variance (KGV), multivariate version of Hoeffding's Φ , Hilbert-Schmidt independence criterion (HSIC), distance covariance, distance correlation, Lancaster three-variable interaction.
- **divergence:** Kullback-Leibler divergence (relative entropy, I directed divergence), Rényi divergence, Tsallis divergence, Sharma-Mittal divergence, Pearson χ^2 divergence (χ^2 distance), Hellinger distance, L_2 divergence, f -divergence (Csiszár-Morimoto divergence, Ali-Silvey distance), maximum mean discrepancy (MMD; kernel distance, current distance), energy distance (N-distance; specifically the Cramer-Von Mises distance), Bhattacharyya distance, non-symmetric Bregman distance (Bregman divergence), symmetric Bregman distance, J-distance (symmetrised Kullback-Leibler divergence, J divergence), K divergence, L divergence, Jensen-Shannon divergence, Jensen-Rényi divergence, Jensen-Tsallis divergence.
- **association measures:** multivariate extensions of Spearman's ρ (Spearman's rank correlation coefficient, grade correlation coefficient), multivariate conditional version of Spearman's ρ , lower and upper tail dependence via conditional Spearman's ρ .
- **cross quantities:** cross-entropy.
- **kernels on distributions:** expected kernel (summation kernel, mean map kernel, set kernel, multi-instance kernel, ensemble kernel; specific convolution kernel), probability product kernel, Bhattacharyya kernel (Bhattacharyya coefficient, Hellinger affinity), Jensen-Shannon kernel, Jensen-Tsallis kernel, exponentiated Jensen-Shannon kernel, exponentiated Jensen-Rényi kernels, exponentiated Jensen-Tsallis kernels.
- **conditional entropy:** conditional Shannon entropy.
- **conditional mutual information:** conditional Shannon mutual information.

- **Web:** <https://bitbucket.org/szzoli/ite-in-python/>. Comments, feedbacks are welcome.

- **Follow ITE:** on

- Bitbucket (<https://bitbucket.org/szzoli/ite-in-python/follow>),
- Twitter (<https://twitter.com/ITEtoolbox>).

- **Mailing list:** <https://groups.google.com/d/forum/itetoolbox>

- **Publications/applications:** Papers using ITE are collected at <https://bitbucket.org/szzoli/ite/wiki>. Feel free to add yours.

- **Author:** Zoltán Szabó (<http://www.cmap.polytechnique.fr/~zoltan.szabo/>).

¹A few nice examples focusing on discrete variables or specialized applications and methods are <http://www.cs.man.ac.uk/~pococka4/MITToolbox.html>, <http://www.cs.tut.fi/~timhome/tim/tim.htm>, <http://cran.r-project.org/web/packages/infotheo>, <http://cran.r-project.org/web/packages/entropy/>, <https://github.com/dit/dit>, <https://pypi.python.org/pypi/universal-divergence/0.2.0>, <https://github.com/baccuslab/shannon>, or <http://fr.mathworks.com/matlabcentral/fileexchange/35625-information-theory-toolbox>.

²See <https://bitbucket.org/szzoli/ite/>. In the sequel we will use 'Matlab ITE' instead of 'Matlab/Octave ITE'.

- **Citing:** If you use the ITE toolbox in your work, please cite it [6].³ The source code also contains references for the individual methods and the quantities estimated.

- **License:** GPLv3(>=).

- **Requirements:**

1. Python 3, SciPy [\ni (typically) NumPy, Matplotlib].⁴ You can get these tools by pip⁵ (\in Python 3 \geq 3.4).

- The system-wide installation is as follows:

```
# python3 -m pip install scipy      # '#' denotes bash prompt (with root rights)
# python3 -m pip install numpy      # if you do not get it by SciPy
# python3 -m pip install matplotlib # -||- (:=same comment)
```

- The user-specific installation is

```
> python3 -m pip install --user scipy # '>' stands for the bash prompt (with normal user)
> python3 -m pip install --user numpy # rights)
> python3 -m pip install --user matplotlib
```

2. Nose, IPython: optional.⁶

Note: Installing Anaconda gives all these tools, with Intel MKL (Math Kernel Library).⁷

The rest of the documentation is structured as follows:

- Section 2 is about the installation of ITE, how to import it and run its built-in demos, and a few usage examples. Section 3 enlists the definitions of the estimated quantities.
- Section A is for developers with details on (i) the directory structure of the toolbox, (ii) how to add new estimators and run doctests, (iii) parameter passing in (certain) meta estimators. Python ITE - Matlab ITE 'correspondence' is the topic of Section B. Section C contains the axiomatic formulation of concordance and dependence.

2 Getting Started: Installation, Built-in Demos, Examples

This section is about the installation and importing of the ITE toolbox, running its built-in demos, followed by a few usage examples.

- Installation: download the ITE archive (https://bitbucket.org/szzoli/ite_in_python/downloads), extract its contents. We will denote the resulting main folder (containing `demos`, `doc`, `ite`, `LICENSE.txt`, ...) as `ite`.
- Start a working session:

```
> ipython3      # see the first bullet point of the note below
>>> import ite  # change first to the ite directory, if it is not on your Python path;
                # '>>>' denotes the prompt in the (I)Python console
```

Note:

- Throughout this documentation for simplicity/efficiency I assume that you use IPython; you might want to do this implicitly via an IDE such as PyCharm⁸.
- You can add the ITE package to the Python path by

³.bib: <http://www.cmap.polytechnique.fr/~zoltan.szabo/ITE.bib>.

⁴See <https://www.python.org/> and <http://www.scipy.org/>.

⁵See <https://pypi.python.org/pypi/pip>.

⁶See <http://nose.readthedocs.io/en/latest/> and <https://ipython.org/>.

⁷See <https://www.continuum.io/downloads>.

⁸See <https://www.jetbrains.com/pycharm/>.

```
>>> import sys
>>> sys.path.insert(1, '/path/to/directory/containing/ite')
```

- Running the built-in demos (see Table 9): Change to the `ite/demos/analytical_values` directory and run the demos. Example:

```
>>> run demo_h_shannon # run ∈ IPython; notice that the '.py' extension could be discarded
```

- Examples: In the first example we estimate $H(\mathbf{y})$, the Shannon entropy [Eq. (1)] of a random variable \mathbf{y} using the k -nearest neighbor method; the estimator is called `BHShannonKnnK` in ITE. \mathbf{y} will be uniformly distributed on the 3-dimensional unit cube ($\in [0, 1]^3$; $d = 3$) from which we have $T = 1000$ samples. The first estimator (`co1` below) relies on the default parameter setting, the second one (`co2`) is based on user-specified parameters. Particularly, in the second case we specify the kNN computation method, the number of neighbors (k) and allow approximation in the kNN phase (`eps`; to speed up computation). For alternative entropy estimators, see Table 1.

Example 1 (Entropy estimation)

```
>>> import ite # import the ITE toolbox (1x)
>>> from numpy.random import rand # we will use 'rand' to create the observations
>>> co1 = ite.cost.BHShannon_KnnK() # initialize the entropy (2nd character = 'H') estimator
>>> print(co1) # print estimator-1
>>> y = rand(1000, 3) # size: number of samples × dimension, {y_t}_{t=1}^{1000}, y_t ∈ ℝ^3
>>> h = co1.estimate(y) # entropy estimation
>>>
>>> co2 = ite.cost.BHShannon_KnnK(knn_method='cKDTree', k=2, eps=0.1) # with other estimator
# parameters
>>> print(co2) # print estimator-2
>>> h2 = co2.estimate(y) # entropy estimation
```

In our second example we consider the estimation of the classical Shannon mutual information [Eq. (6)]. The random variable (\mathbf{y}) of interest is partitioned into 3 blocks: $\mathbf{y} = [\mathbf{y}^1; \mathbf{y}^2; \mathbf{y}^3]$ ($\mathbf{y}^1 \in \mathbb{R}^2$, $\mathbf{y}^2 \in \mathbb{R}^3$, $\mathbf{y}^3 \in \mathbb{R}^4$) and we want to get $I(\mathbf{y}^1, \mathbf{y}^2, \mathbf{y}^3)$, the mutual information of \mathbf{y}^m -s. $T = 2000$ samples are used for estimation. We estimate the mutual information of \mathbf{y}^m -s from Kullback-Leibler divergence [see Eq. (6) and Eq. (24)]. These type of *derived* estimators are called meta estimators in ITE (`MIShannon_DKL`; 1st character = 'M'). 'Base' estimators refer to non-derived ones, as it was the case in Example 1 (`BHShannon_KnnK`: 1st character = 'B').

Example 2 (Mutual information estimation (association measure: similarly))

```
>>> from numpy.random import randn # we will use 'randn' to create the observations
>>> from numpy import array # an 'array' will contain the subspace dimensions: [d1; d2; d3]
>>> co = ite.cost.MIShannon_DKL() # initialize the mutual information estimator
# (MIShannon_DKL: 2nd character = 'I')
>>> ds = array([2, 3, 4]) # y^1 ∈ ℝ^2, y^2 ∈ ℝ^3, y^3 ∈ ℝ^4, d = d1 + d2 + d3 = 2 + 3 + 4 = 9
>>> t = 2000 # number of samples
>>> y = randn(t, sum(ds)) # size: number of samples × dimension
>>> i = co.estimate(y, ds) # estimate mutual information
```

Alternative mutual information measures/techniques are listed in Table 2. Estimation of association measures is analogous, the available methods are covered in Table 4.

In our third example we estimate $D(\mathbf{y}^1, \mathbf{y}^2)$, the Kullback-Leibler divergence between two random quantities \mathbf{y}^1 and \mathbf{y}^2 [see Eq. (24)⁹], via k -nearest neighbors (`BDKL_KnnK`). We have $T_1 = 2000$ samples from \mathbf{y}^1 and $T_2 = 3000$ samples from \mathbf{y}^2 .

⁹We identify random variables with their distributions or their pdf-s (probability density functions; meant w.r.t. the Lebesgue measure). Many of the information theoretical quantities can be formulated more generally, but to keep the presentation simple we will avoid going into measure theoretical details/technicalities.

Example 3 (Divergence estimation (cross quantity, kernel on distributions: analogously))

```
>>> from numpy.random import randn # 'randn' is used to generate our observations
>>> co = ite.cost.BDKL_KnnK() # initialize the divergence (2nd character = 'D') estimator
>>> dim = 3 #  $\mathbf{y}^1 \in \mathbb{R}^3$ ,  $\mathbf{y}^2 \in \mathbb{R}^3$ 
>>> t1, t2 = 2000, 3000 # number of samples from  $\mathbf{y}^1$  and  $\mathbf{y}^2$ 
>>> y1 = randn(t1, dim) # size: number of samples1  $\times$  dimension,  $\{\mathbf{y}_t^1\}$ 
>>> y2 = randn(t2, dim) # size: number of samples2  $\times$  dimension,  $\{\mathbf{y}_t^2\}$ 
>>> d = co.estimate(y1, y2) # estimate KL divergence
```

For other divergence measures or estimators see Table 3. The estimation of cross quantities and kernels on distributions (Table 5, Table 6) can be carried out in the same way.

In our fourth example, we focus on the estimation of the conditional Shannon entropy of \mathbf{y}^1 given \mathbf{y}^2 [see Eq. (68)]. We have $T = 5000$ samples from the joint distribution of $\mathbf{y} = [\mathbf{y}^1, \mathbf{y}^2]$; it is assumed to be Gaussian below.

Example 4 (Conditional entropy estimation)

```
>>> from numpy import dot # create observations
>>> from numpy.random import rand, multivariate_normal # -||-
>>> dim1, dim2 = 1, 2 #  $\mathbf{y}^1 \in \mathbb{R}^1$ ,  $\mathbf{y}^2 \in \mathbb{R}^2$ 
>>> dim = dim1 + dim2 #  $\mathbf{y} = [\mathbf{y}^1, \mathbf{y}^2] \in \mathbb{R}^{1+2=3}$ 
>>> t = 5000 # number of samples
>>> co = ite.cost.BcondHShannon_HShannon() # initialize the conditional entropy ('condH')
# estimator
>>> m, l = rand(dim), rand(dim, dim) # mean ( $\mathbf{m}$ )
>>> c = dot(l, l.T) # covariance ( $\Sigma$ ),  $\mathbf{y} = N(\mathbf{m}, \Sigma)$ 
>>> y = multivariate_normal(m, c, t) #  $\{\mathbf{y}_t\}_{t=1}^{5000}$ ,  $\mathbf{y}_t = [\mathbf{y}_t^1, \mathbf{y}_t^2] \in \mathbb{R}^3$ 
>>> cond_h = co.estimate(y, dim1) # estimate conditional entropy
```

In our fifth example, the task is to estimate the conditional Shannon mutual information of \mathbf{y}^1 and \mathbf{y}^2 given \mathbf{y}^3 [see Eq. (69)]. We are given $T = 3000$ samples from the joint distribution of $\mathbf{y} = [\mathbf{y}^1; \mathbf{y}^2; \mathbf{y}^3]$; in the example below it is Gaussian.

Example 5 (Conditional mutual information estimation)

```
>>> from numpy import dot, array # create observations
>>> from numpy.random import rand, multivariate_normal # -||-
>>> ds = array([1, 2, 3]) #  $\mathbf{y}^1 \in \mathbb{R}$ ,  $\mathbf{y}^2 \in \mathbb{R}^2$ ,  $\mathbf{y}^3 \in \mathbb{R}^3$ 
>>> dim = sum(ds) #  $d = d_1 + d_2 + d_3 = 1 + 2 + 3 = 6$ ,
#  $\mathbf{y} = [\mathbf{y}^1, \mathbf{y}^2, \mathbf{y}^3] \in \mathbb{R}^6$ 
>>> t = 3000 # number of samples
>>> co = ite.cost.BcondIShannon_HShannon() # initialize the conditional mutual information
# ('condI') estimator
>>> m, l = rand(dim), rand(dim, dim) # mean ( $\mathbf{m}$ )
>>> c = dot(l, l.T) # covariance ( $\Sigma$ ),  $\mathbf{y} = N(\mathbf{m}, \Sigma)$ 
>>> y = multivariate_normal(m, c, t) #  $\{\mathbf{y}_t\}_{t=1}^{3000}$ ,  $\mathbf{y}_t = [\mathbf{y}_t^1, \mathbf{y}_t^2, \mathbf{y}_t^3]$ 
>>> cond_i = co.estimate(y, ds) # estimate conditional mutual information
```

3 Estimated Quantities and Estimators

This section provides the definitions of the available information theoretical quantities in ITE: Section 3.1 focuses on unconditional quantities, Section 3.2 contains the conditional ones. Section 3.3 is about the estimators of these quantities.

3.1 Unconditional Quantities

This part is structured as follows: entropy (Section 3.1.1), mutual information (Section 3.1.2), divergence (Section 3.1.3), association measure (Section 3.1.4) cross quantity (Section 3.1.5), kernel on distributions (Section 3.1.6).

3.1.1 Entropy

- Notation: $\mathbb{R}^d \ni \mathbf{y} \sim f$, in other words the d -dimensional random variable \mathbf{y} has density f .
- Goal: We want to estimate the entropy of $\mathbf{y} \in \mathbb{R}^d$ from which we have i.i.d. (independent identically distributed) samples, $\{\mathbf{y}_t\}_{t=1}^T$ ($\mathbf{y}_t \in \mathbb{R}^d$, $t = 1, \dots, T$).
- Definition: The Shannon entropy (H), Rényi entropy ($H_{R,\alpha}$), Tsallis entropy ($H_{T,\alpha}$; also called Havrda and Charvát entropy), Sharma-Mittal entropy ($H_{SM,\alpha,\beta}$), Φ -entropy ($H_{\Phi,w}$; f -entropy¹⁰) are defined as¹¹

$$H(\mathbf{y}) = - \int_{\mathbb{R}^d} f(\mathbf{u}) \log f(\mathbf{u}) d\mathbf{u}, \quad (1)$$

$$H_{R,\alpha}(\mathbf{y}) = \frac{1}{1-\alpha} \log \int_{\mathbb{R}^d} f^\alpha(\mathbf{u}) d\mathbf{u}, \quad (\alpha \neq 1) \quad \lim_{\alpha \rightarrow 1} H_{R,\alpha} = H, \quad (2)$$

$$H_{T,\alpha}(\mathbf{y}) = \frac{1}{\alpha-1} \left[1 - \int_{\mathbb{R}^d} f^\alpha(\mathbf{u}) d\mathbf{u} \right] = \frac{e^{(1-\alpha)H_{R,\alpha}(\mathbf{y})} - 1}{1-\alpha}, \quad (\alpha \neq 1) \quad \lim_{\alpha \rightarrow 1} H_{T,\alpha} = H, \quad (3)$$

$$H_{SM,\alpha,\beta}(\mathbf{y}) = \frac{1}{1-\beta} \left[\left(\int_{\mathbb{R}^d} f^\alpha(\mathbf{u}) d\mathbf{u} \right)^{\frac{1-\beta}{1-\alpha}} - 1 \right], \quad (\alpha > 0, \alpha \neq 1, \beta \neq 1), \quad (4)$$

$$H_{\Phi,w}(y) = \int_{\mathbb{R}} f(u) \Phi(f(u)) w(u) du, \quad (\mathbb{R} \ni y \sim f). \quad (5)$$

- Note:

$$- H_{SM,\alpha,\beta}: \lim_{\beta \rightarrow 1} H_{SM,\alpha,\beta}(\mathbf{y}) = H_{R,\alpha}(\mathbf{y}), \quad H_{SM,\alpha,\alpha}(\mathbf{y}) = H_{T,\alpha}(\mathbf{y}), \quad \lim_{(\alpha,\beta) \rightarrow (1,1)} H_{SM,\alpha,\beta}(\mathbf{y}) = H(\mathbf{y}).$$

3.1.2 Mutual Information

- Notation: $\mathbb{R}^d \ni \mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M] \sim f$, $\mathbb{R}^{d_m} \ni \mathbf{y}^m \sim f_m$ ($d = \sum_{m=1}^M d_m$). f_S ($S \subseteq \{1, \dots, M\}$) stands for the associated marginals; for example $f_{\{1,2\}}$ is the density function of $[\mathbf{y}^1; \mathbf{y}^2]$. $\widehat{\mathbf{y}}^m$ denotes an identically distributed copy of \mathbf{y}^m . ‘ \perp ’ means independence, ‘ \vee ’ denotes the logical ‘or’. \mathbb{E} is for expectation, *cov* is covariance, *var* denotes variance, $i = \sqrt{-1}$, $\langle \cdot, \cdot \rangle_2$ is the Euclidean inner product, $\Gamma(t) = \int_0^\infty x^{t-1} e^{-x} dx$ is the Gamma function. Let us define the weight function $w(\mathbf{u}^1, \mathbf{u}^2) = \frac{1}{c(d_1,\alpha)c(d_2,\alpha)[\|\mathbf{u}^1\|_2]^{d_1+\alpha}[\|\mathbf{u}^2\|_2]^{d_2+\alpha}}$ with $\alpha \in (0, 2)$, where $c(d, \alpha) = \frac{2\pi^{\frac{d}{2}} \Gamma(1-\frac{\alpha}{2})}{\alpha 2^\alpha \Gamma(\frac{d+\alpha}{2})}$. $\varphi_{12}(\mathbf{u}^1, \mathbf{u}^2) = \mathbb{E}_{\mathbf{y}^1 \mathbf{y}^2} \left[e^{i\langle \mathbf{u}^1, \mathbf{y}^1 \rangle + i\langle \mathbf{u}^2, \mathbf{y}^2 \rangle} \right]$, $\varphi_j(\mathbf{u}^j) = \mathbb{E}_{\mathbf{y}^j} \left[e^{i\langle \mathbf{u}^j, \mathbf{y}^j \rangle} \right]$, ($j = 1, 2$) are the characteristic functions of $[\mathbf{y}^1; \mathbf{y}^2]$, \mathbf{y}^1 and \mathbf{y}^2 . Given a reproducing kernel k , $\mathcal{H}(k)$ is the associated RKHS (reproducing kernel Hilbert space). $\otimes_{i=1}^3 \mathcal{H}(k_i)$ denotes the tensor product of $\mathcal{H}(k_i)$ -s. $\mu_q = \int k(\cdot, u) q(u) du = \mathbb{E}_{u \sim q} [k(\cdot, u)]$ is the mean embedding of q to $\mathcal{H}(k)$; q is often a pdf [see Eq. (20), (21) where q is not a pdf, it can be negative]. F : cdf (cumulative density function) of $\mathbf{y} = [y^1; \dots; y^d]$, F_i : cdf of y^i . C : copula of \mathbf{y} , i.e. $F(\mathbf{y}) = C(F_1(y^1), \dots, F_d(y^d))$, in other words $C(\mathbf{u}) = \mathbb{P}(\mathbf{U} \leq \mathbf{u})$ where $\mathbf{U} = [F_1(y^1); \dots; F_d(y^d)] \in [0, 1]^d$. $\Pi(u_1, \dots, u_d) = \prod_{i=1}^d u_i$ is the product copula. f_U : uniform density on $[0, 1]^M$.
- Goal: We consider the estimation of the mutual information of the d_m -dimensional components (\mathbf{y}^m) of the random variable \mathbf{y} using an i.i.d. sample set $\{\mathbf{y}_t\}_{t=1}^T$ from \mathbf{y} .
- Definition: The Shannon mutual information (I ; also known as total correlation or multi-information), Rényi mutual information ($I_{R,\alpha}$), Tsallis mutual information ($I_{T,\alpha}$), χ^2 mutual information (I_{χ^2} ; for $M = 2$ also called squared-loss mutual information; mean square contingency $= \sqrt{I_{\chi^2}}$), L_2 mutual information (I_{L_2}), copula-based kernel dependency (I_c), kernel canonical correlation analysis (I_{KCCA} ; $KCCA$), kernel generalized variance (I_{KGV} ; KGV), multivariate version of Hoeffding’s Φ (I_Φ), Hilbert-Schmidt independence criterion (I_{HSIC} ; $HSIC$), distance covariance (I_{dCov}), distance correlation (I_{dCor}), Lancaster 3-variable interaction (I_{3-Lanc}), three-variable joint independence measure ($I_{3-joint}$) are defined as

$$I(\mathbf{y}^1, \dots, \mathbf{y}^M) = \int_{\mathbb{R}^d} f(\mathbf{u}^1, \dots, \mathbf{u}^M) \log \left[\frac{f(\mathbf{u}^1, \dots, \mathbf{u}^M)}{\prod_{m=1}^M f_m(\mathbf{u}^m)} \right] d\mathbf{u}^1 \dots d\mathbf{u}^M, \quad (6)$$

¹⁰Since f also denotes the density in Eq. (5), we refer to the quantity as the Φ -entropy.

¹¹Here and in the sequel \log denotes natural logarithm, i.e., the unit of the information theoretical measures is nat.

$$= D \left(f, \prod_{m=1}^M f_m \right) \quad [\text{see (24) for the definition of } D], \quad (7)$$

$$I_{R,\alpha}(\mathbf{y}^1, \dots, \mathbf{y}^M) = D_{R,\alpha} \left(f, \prod_{m=1}^M f_m \right), \text{ for } D_{R,\alpha}, \text{ see (25)}, \quad (8)$$

$$I_{T,\alpha}(\mathbf{y}^1, \dots, \mathbf{y}^M) = D_{T,\alpha} \left(f, \prod_{m=1}^M f_m \right), \text{ for } D_{T,\alpha}, \text{ see (26)}, \quad (9)$$

$$I_{\chi^2}(\mathbf{y}^1, \dots, \mathbf{y}^M) = D_{\chi^2} \left(f, \prod_{m=1}^M f_m \right), \text{ for } D_{\chi^2}, \text{ see (28)}, \quad (10)$$

$$I_{L_2}(\mathbf{y}^1, \dots, \mathbf{y}^M) = D_{L_2} \left(f, \prod_{m=1}^M f_m \right), \text{ for } D_{L_2}, \text{ see (30)}, \quad (11)$$

$$I_c(\mathbf{y}^1, \dots, \mathbf{y}^M) = D_{\text{MMD}}(f_{\mathbf{z}}, f_U), \text{ for } D_{\text{MMD}} \text{ see (32)}, \quad \mathbf{z} = [F_1(\mathbf{y}^1); \dots; F_M(\mathbf{y}^M)] \in \mathbb{R}^M, \quad (12)$$

$$I_{\text{KCCA}}(\mathbf{y}^1, \mathbf{y}^2) = \sup_{g_1 \in \mathcal{H}(k_1), g_2 \in \mathcal{H}(k_2)} \frac{\text{cov}[g_1(\mathbf{y}^1), g_2(\mathbf{y}^2)]}{\sqrt{\text{var}[g_1(\mathbf{y}^1)] + \kappa \|g_1\|_{\mathcal{H}(k_1)}^2} \sqrt{\text{var}[g_2(\mathbf{y}^2)] + \kappa \|g_2\|_{\mathcal{H}(k_2)}^2}}, \quad (\kappa > 0), \quad (13)$$

$$I_{\text{KGV}}(\mathbf{y}^1, \dots, \mathbf{y}^M) = -\frac{1}{2} \log \left[\frac{\det(\mathbf{C})}{\prod_{m=1}^M \det(\mathbf{C}^{m,m})} \right], \quad \mathbf{C} = [\mathbf{C}^{i,j}], \quad \varphi(\mathbf{y}) := [\varphi_m(\mathbf{y}^m)]_{m=1}^M, \quad (14)$$

$$\mathbf{C}^{i,j} = \text{cov}[\varphi_i(\mathbf{y}^i), \varphi_j(\mathbf{y}^j)], \quad (15)$$

$$I_{\Phi}(\mathbf{y}^1, \dots, \mathbf{y}^d) = I_{\Phi}(C) = \left(h_2(d) \int_{[0,1]^d} [C(\mathbf{u}) - \Pi(\mathbf{u})]^2 d\mathbf{u} \right)^{\frac{1}{2}}, \quad (16)$$

$$h_2(d) = \left(\frac{2}{(d+1)(d+2)} - \frac{1}{2^d} \frac{d!}{\prod_{i=0}^d (i + \frac{1}{2})} + \frac{1}{3^d} \right)^{-1},$$

$$I_{\text{HSIC}}(\mathbf{y}^1, \mathbf{y}^2) = \|C_{\mathbf{y}^1 \mathbf{y}^2}\|_{\text{HS}}^2, \quad C_{\mathbf{y}^1 \mathbf{y}^2} = \mathbb{E}_{\mathbf{y}^1 \mathbf{y}^2} ([k_1(\cdot, \mathbf{y}^1) - \boldsymbol{\mu}_1] \otimes [k_2(\cdot, \mathbf{y}^2) - \boldsymbol{\mu}_2]), \quad (17)$$

$$I_{\text{dCov}}(\mathbf{y}^1, \mathbf{y}^2) = \|\varphi_{12} - \varphi_1 \varphi_2\|_{L_w^2} = \sqrt{\int_{\mathbb{R}^{d_1+d_2}} |\varphi_{12}(\mathbf{u}^1, \mathbf{u}^2) - \varphi_1(\mathbf{u}^1) \varphi_2(\mathbf{u}^2)|^2 w(\mathbf{u}^1, \mathbf{u}^2) d\mathbf{u}^1 d\mathbf{u}^2}, \quad (18)$$

$$I_{\text{dCor}}(\mathbf{y}^1, \mathbf{y}^2) = \begin{cases} \frac{I_{\text{dCov}}(\mathbf{y}^1, \mathbf{y}^2)}{\sqrt{I_{\text{dVar}}(\mathbf{y}^1, \mathbf{y}^1) I_{\text{dVar}}(\mathbf{y}^2, \mathbf{y}^2)}}, & \text{if } I_{\text{dVar}}(\mathbf{y}^1, \mathbf{y}^1) I_{\text{dVar}}(\mathbf{y}^2, \mathbf{y}^2) > 0, \\ 0, & \text{otherwise,} \end{cases} \quad (19)$$

$$I_{\text{dVar}}(\mathbf{y}^j, \mathbf{y}^j) = \|\varphi_{jj} - \varphi_j \varphi_j\|_{L_w^2},$$

$$I_{3\text{-Lanc}}(\mathbf{y}^1, \mathbf{y}^2, \mathbf{y}^3) = \|\mu_{L(f)}\|_{\otimes_{i=1}^3 \mathcal{H}(k_i)}^2, \quad L(f) = f - f_{12}f_3 - f_{23}f_1 - f_{13}f_2 + 2f_1f_2f_3. \quad (20)$$

$$I_{3\text{-joint}}(\mathbf{y}^1, \mathbf{y}^2, \mathbf{y}^3) = \|\mu_{J(f)}\|_{\otimes_{i=1}^3 \mathcal{H}(k_i)}^2, \quad J(f) = f - f_1f_2f_3, \quad (21)$$

where $C_{\mathbf{y}^1 \mathbf{y}^2}$ is the so-called cross-covariance operator, $\|\cdot\|_{\text{HS}}$ is the Hilbert-Schmidt norm, $L(f)$ is the Lancaster interaction measure, μ is the mean embedding to either $\mathcal{H}(k_i)$ [for I_{HSIC}] or to $\otimes_{i=1}^3 \mathcal{H}(k_i)$ [in case of $I_{3\text{-Lanc}}$ and $I_{3\text{-joint}}$], φ_m is the canonical feature map associated to kernel k_m .

• Note:

1. $I(\mathbf{y}^1, \dots, \mathbf{y}^M) \geq 0$, $I(\mathbf{y}^1, \dots, \mathbf{y}^M) = 0 \Leftrightarrow \mathbf{y}^m$ -s are jointly independent.
2. $I_{R,\alpha}, I_{T,\alpha}$: $\lim_{\alpha \rightarrow 1} I_{R,\alpha}(\mathbf{y}) = \lim_{\alpha \rightarrow 1} I_{T,\alpha}(\mathbf{y}) = I(\mathbf{y})$.
3. I_{KCCA} : KCCA captures the maximal correlation in the $\mathcal{H}(k_i)$ feature spaces. It can be generalized to $M \geq 2$ components to measure pairwise independence; this extension is available in ITE.
4. I_{KGV} : KGV is the extension of the analytical expression of mutual information holding for Gaussian variables (\mathbf{y}).
5. I_{Φ} : It is a multivariate ($M \geq 2$) extension of Hoeffding's Φ capturing the deviation from the product copula in $L^2([0, 1]^d)$ sense. $h_2(d)$ is a normalizing constant ensuring that $I_{\Phi}(C) \in [0, 1]$ for any copula C .

6. I_{HSIC} : It can also be extended to the $M \geq 2$ case to measure pairwise independence; it is in ITE.

7. I_{dCov} :

- It measures independence by the difference of the joint characteristic function and the product of the marginals, in L_w^2 sense.
- For $\alpha = 1$ the distance covariance can be rewritten in terms of pairwise distances

$$\begin{aligned} I_{\text{dCov}}(\mathbf{y}^1, \mathbf{y}^2) &= \mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} \mathbb{E}_{\widetilde{\mathbf{y}}^1, \widetilde{\mathbf{y}}^2} \left[\left\| \mathbf{y}^1 - \widetilde{\mathbf{y}}^1 \right\|_2 \left\| \mathbf{y}^2 - \widetilde{\mathbf{y}}^2 \right\|_2 \right] + \mathbb{E}_{\mathbf{y}^1, \widetilde{\mathbf{y}}^1} \left[\left\| \mathbf{y}^1 - \widetilde{\mathbf{y}}^1 \right\|_2 \right] \mathbb{E}_{\mathbf{y}^2, \widetilde{\mathbf{y}}^2} \left[\left\| \mathbf{y}^2 - \widetilde{\mathbf{y}}^2 \right\|_2 \right] \\ &\quad - 2 \mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} \left[\mathbb{E}_{\widetilde{\mathbf{y}}^1} \left\| \mathbf{y}^1 - \widetilde{\mathbf{y}}^1 \right\|_2 \mathbb{E}_{\widetilde{\mathbf{y}}^2} \left\| \mathbf{y}^2 - \widetilde{\mathbf{y}}^2 \right\|_2 \right]. \end{aligned} \quad (22)$$

This form has a natural extension to semimetric spaces $[\mathbf{y}^1 \in (\mathcal{Y}_1, \rho_1), \mathbf{y}^2 \in (\mathcal{Y}_2, \rho_2)]$ of negative type:

$$\begin{aligned} I_{\text{dCov}}(\mathbf{y}^1, \mathbf{y}^2) &= \mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} \mathbb{E}_{\widetilde{\mathbf{y}}^1, \widetilde{\mathbf{y}}^2} \left[\rho_1(\mathbf{y}^1, \widetilde{\mathbf{y}}^1) \rho_2(\mathbf{y}^2, \widetilde{\mathbf{y}}^2) \right] + \mathbb{E}_{\mathbf{y}^1, \widetilde{\mathbf{y}}^1} \left[\rho_1(\mathbf{y}^1, \widetilde{\mathbf{y}}^1) \right] \mathbb{E}_{\mathbf{y}^2, \widetilde{\mathbf{y}}^2} \left[\rho_2(\mathbf{y}^2, \widetilde{\mathbf{y}}^2) \right] \\ &\quad - 2 \mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} \left(\mathbb{E}_{\widetilde{\mathbf{y}}^1} \left[\rho_1(\mathbf{y}^1, \widetilde{\mathbf{y}}^1) \right] \mathbb{E}_{\widetilde{\mathbf{y}}^2} \left[\rho_2(\mathbf{y}^2, \widetilde{\mathbf{y}}^2) \right] \right), \end{aligned}$$

which is proportional to HSIC (determined by kernel k):

$$\begin{aligned} I_{\text{dCov}}(\mathbf{y}^1, \mathbf{y}^2) &= 2I_{\text{HSIC}}(\mathbf{y}^1, \mathbf{y}^2), \quad k((\mathbf{u}_1, \mathbf{v}_1), (\mathbf{u}_2, \mathbf{v}_2)) = k_1(\mathbf{u}_1, \mathbf{u}_2)k_2(\mathbf{v}_1, \mathbf{v}_2), \\ &\quad \rho_i(\mathbf{u}, \mathbf{v}) = k_i(\mathbf{u}, \mathbf{u}) + k_i(\mathbf{v}, \mathbf{v}) - 2k_i(\mathbf{u}, \mathbf{v}). \end{aligned} \quad (23)$$

8. I_{dCor} : is the normalized variant of I_{dCov} ; $I_{\text{dCor}}(\mathbf{y}^1, \mathbf{y}^2) \in [0, 1]$. It is zero iff \mathbf{y}^1 and \mathbf{y}^2 are independent.

9. $I_{3\text{-Lanc}}$: guaranteed to be zero, if f can be factorised as a product of its (possible) multidimensional marginals, i.e. $([\mathbf{y}^1; \mathbf{y}^2] \perp\!\!\!\perp \mathbf{y}^3) \vee ([\mathbf{y}^1; \mathbf{y}^3] \perp\!\!\!\perp \mathbf{y}^2) \vee ([\mathbf{y}^2; \mathbf{y}^3] \perp\!\!\!\perp \mathbf{y}^1) \Rightarrow L(f) = 0$. For example, $[\mathbf{y}^1; \mathbf{y}^2] \perp\!\!\!\perp \mathbf{y}^3$ stands for $f = f_{12}f_3$.

10. $I_{3\text{-joint}}$: measures joint independence in $\otimes_{i=1}^3 \mathcal{H}(k_i)$.

3.1.3 Divergence

- Notation: $\mathbb{R}^d \ni \mathbf{y}_1 \sim f_1, \mathbb{R}^d \ni \mathbf{y}_2 \sim f_2$. $\pi_1 \widetilde{\mathbf{y}}^1 + \pi_2 \mathbf{y}^2$ is the mixture distribution obtained from \mathbf{y}^1 and \mathbf{y}^2 with π_1, π_2 weights ($\pi_1, \pi_2 > 0, \pi_1 + \pi_2 = 1$). $\widetilde{\mathbf{y}}^1$ and \mathbf{y}^2 are identically distributed copies of \mathbf{y}^1 and \mathbf{y}^2 . $\text{supp}(f_i)$ is the support of the pdf f_i .
- Goal: Given independent, i.i.d. samples from f_1 and f_2 , $\{\mathbf{y}_t^1\}_{t=1}^{T_1}$ and $\{\mathbf{y}_t^2\}_{t=1}^{T_2}$, we want to estimate the divergence of the two underlying random variables (\mathbf{y}_1 and \mathbf{y}_2).
- Definition: The Kullback-Leibler divergence (D ; also called relative entropy or I directed divergence), Rényi divergence ($D_{\text{R},\alpha}$), Tsallis divergence ($D_{\text{T},\alpha}$), Sharma-Mittal divergence ($D_{\text{SM},\alpha,\beta}$), Pearson χ^2 divergence (D_{χ^2} ; also called χ^2 distance), Hellinger distance (D_{H}), L_2 divergence (D_{L_2}), (Csiszár) f-divergence (D_f ; also called Csiszár-Morimoto divergence or Ali-Silvey distance), maximum mean discrepancy (D_{MMD} ; MMD, also called kernel distance, current distance), energy distance (D_{EnDist} ; also called N-distance), Bhattacharyya distance (D_{B}), non-symmetric Bregman distance ($D_{\text{NB},\alpha}$; also called Bregman divergence), symmetric Bregman distance ($D_{\text{SB},\alpha}$), J-distance (D_{J} ; symmetrised Kullback-Leibler divergence, J divergence), K divergence (D_{K}), L divergence (D_{L}) Jensen-Shannon divergence (D_{JS}^π), Jensen-Rényi divergence ($D_{\text{JR},\alpha}^\pi$), Jensen-Tsallis divergence ($D_{\text{JT},\alpha}$) are defined as

$$D(f_1, f_2) = \int_{\mathbb{R}^d} f_1(\mathbf{u}) \log \left[\frac{f_1(\mathbf{u})}{f_2(\mathbf{u})} \right] d\mathbf{u}, \quad (24)$$

$$D_{\text{R},\alpha}(f_1, f_2) = \frac{1}{\alpha - 1} \log \int_{\mathbb{R}^d} f_1^\alpha(\mathbf{u}) f_2^{1-\alpha}(\mathbf{u}) d\mathbf{u}, \quad (\alpha \in \mathbb{R} \setminus \{1\}), \quad (25)$$

$$D_{\text{T},\alpha}(f_1, f_2) = \frac{1}{\alpha - 1} \left(\int_{\mathbb{R}^d} f_1^\alpha(\mathbf{u}) f_2^{1-\alpha}(\mathbf{u}) d\mathbf{u} - 1 \right), \quad (\alpha \in \mathbb{R} \setminus \{1\}), \quad (26)$$

$$D_{\text{SM},\alpha,\beta}(f_1, f_2) = \frac{1}{\beta - 1} \left[\left(\int_{\mathbb{R}^d} [f_1(\mathbf{u})]^\alpha [f_2(\mathbf{u})]^{1-\alpha} d\mathbf{u} \right)^{\frac{1-\beta}{1-\alpha}} - 1 \right], \quad (\alpha \neq 1, \beta \neq 1), \quad (27)$$

$$D_{\chi^2}(f_1, f_2) = \int_{\text{supp}(f_1) \cup \text{supp}(f_2)} \frac{[f_1(\mathbf{u}) - f_2(\mathbf{u})]^2}{f_2(\mathbf{u})} d\mathbf{u} = \int_{\text{supp}(f_1) \cup \text{supp}(f_2)} \frac{[f_1(\mathbf{u})]^2}{f_2(\mathbf{u})} d\mathbf{u} - 1, \quad (28)$$

$$D_{\text{H}}(f_1, f_2) = \sqrt{\frac{1}{2} \int_{\mathbb{R}^d} [\sqrt{f_1(\mathbf{u})} - \sqrt{f_2(\mathbf{u})}]^2 d\mathbf{u}} = \sqrt{1 - \int_{\mathbb{R}^d} \sqrt{f_1(\mathbf{u})} \sqrt{f_2(\mathbf{u})} d\mathbf{u}}, \quad (29)$$

$$D_{\text{L}_2}(f_1, f_2) = \sqrt{\int_{\mathbb{R}^d} [f_1(\mathbf{u}) - f_2(\mathbf{u})]^2 d\mathbf{u}}, \quad (30)$$

$$D_f(f_1, f_2) = \int_{\mathbb{R}^d} f \left[\frac{f_1(\mathbf{u})}{f_2(\mathbf{u})} \right] f_2(\mathbf{u}) d\mathbf{u}, \quad f: \text{convex}, f(1) = 0, \quad (31)$$

$$D_{\text{MMD}}(f_1, f_2) = \|\mu_1 - \mu_2\|_{\mathcal{H}(k)}, \quad (32)$$

$$\begin{aligned} D_{\text{EnDist}}(f_1, f_2) &= 2\mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} [\rho(\mathbf{y}^1, \mathbf{y}^2)] - \mathbb{E}_{\mathbf{y}^1, \widetilde{\mathbf{y}}^1} [\rho(\mathbf{y}^1, \widetilde{\mathbf{y}}^1)] - \mathbb{E}_{\mathbf{y}^2, \widetilde{\mathbf{y}}^2} [\rho(\mathbf{y}^2, \widetilde{\mathbf{y}}^2)] \xrightarrow{\text{specifically: } \rho(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_2^\alpha} \\ &= 2\mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} \|\mathbf{y}^1 - \mathbf{y}^2\|_2^\alpha - \mathbb{E}_{\mathbf{y}^1, \widetilde{\mathbf{y}}^1} \|\mathbf{y}^1 - \widetilde{\mathbf{y}}^1\|_2^\alpha - \mathbb{E}_{\mathbf{y}^2, \widetilde{\mathbf{y}}^2} \|\mathbf{y}^2 - \widetilde{\mathbf{y}}^2\|_2^\alpha, \quad \alpha \in (0, 2), \end{aligned} \quad (33)$$

$$D_{\text{B}}(f_1, f_2) = -\log \left(\int_{\mathbb{R}^d} \sqrt{f_1(\mathbf{u})} \sqrt{f_2(\mathbf{u})} d\mathbf{u} \right), \quad (34)$$

$$D_{\text{NB}, \alpha}(f_1, f_2) = \int_{\mathbb{R}^d} \left[f_2^\alpha(\mathbf{u}) + \frac{1}{\alpha - 1} f_1^\alpha(\mathbf{u}) - \frac{\alpha}{\alpha - 1} f_1(\mathbf{u}) f_2^{\alpha-1}(\mathbf{u}) \right] d\mathbf{u}, \quad (\alpha \neq 1), \quad (35)$$

$$D_{\text{SB}, \alpha}(f_1, f_2) = \frac{1}{\alpha} [D_{\text{NB}, \alpha}(f_1, f_2) + D_{\text{NB}, \alpha}(f_2, f_1)], \quad (\alpha \neq 1) \quad (36)$$

$$= \frac{1}{\alpha - 1} \int_{\mathbb{R}^d} f_1^\alpha(\mathbf{u}) + f_2^\alpha(\mathbf{u}) - f_1(\mathbf{u}) f_2^{\alpha-1}(\mathbf{u}) - f_2(\mathbf{u}) f_1^{\alpha-1}(\mathbf{u}) d\mathbf{u}, \quad (37)$$

$$D_{\text{J}}(f_1, f_2) = D(f_1, f_2) + D(f_2, f_1), \quad (38)$$

$$D_{\text{K}}(f_1, f_2) = D \left(f_1, \frac{f_1 + f_2}{2} \right), \quad (39)$$

$$D_{\text{L}}(f_1, f_2) = D_{\text{K}}(f_1, f_2) + D_{\text{K}}(f_2, f_1), \quad (40)$$

$$D_{\text{JS}}^\pi(f_1, f_2) = H(\pi_1 \mathbf{y}^1 + \pi_2 \mathbf{y}^2) - [\pi_1 H(\mathbf{y}^1) + \pi_2 H(\mathbf{y}^2)], \quad (41)$$

$$D_{\text{JR}, \alpha}^\pi(f_1, f_2) = H_{\text{R}, \alpha}(\pi_1 \mathbf{y}^1 + \pi_2 \mathbf{y}^2) - [\pi_1 H_{\text{R}, \alpha}(\mathbf{y}^1) + \pi_2 H_{\text{R}, \alpha}(\mathbf{y}^2)], \quad (0 < \alpha \neq 1), \quad (42)$$

$$D_{\text{JT}, \alpha}(f_1, f_2) = H_{\text{T}, \alpha} \left(\frac{\mathbf{y}^1 + \mathbf{y}^2}{2} \right) - \frac{H_{\text{T}, \alpha}(\mathbf{y}^1) + H_{\text{T}, \alpha}(\mathbf{y}^2)}{2}, \quad (\alpha \neq 1), \quad (43)$$

where $\mu_m = \mathbb{E}_{\mathbf{y}_m \sim f_m} [k(\cdot, \mathbf{y}_m)]$ is the mean embedding of f_m to the $\mathcal{H}(k)$ RKHS.

• Note:

– D : $D(f_1, f_2) \geq 0$. $D(f_1, f_2) = 0 \Leftrightarrow f_1 = f_2$. It is a specific f-divergence [see (31)] with $f(t) = t \log(t)$.

– $D_{\text{R}, \alpha}$, $D_{\text{T}, \alpha}$:

* $\lim_{\alpha \rightarrow 1} D_{\text{R}, \alpha}(f_1, f_2) = \lim_{\alpha \rightarrow 1} D_{\text{T}, \alpha}(f_1, f_2) = D(f_1, f_2)$.

* $\alpha < 0 \Rightarrow D_{\text{R}, \alpha}(f_1, f_2) \leq 0, D_{\text{T}, \alpha}(f_1, f_2) \leq 0$.

* $\alpha = 0 \Rightarrow D_{\text{R}, \alpha}(f_1, f_2) = D_{\text{T}, \alpha}(f_1, f_2) = 0$.

* $\alpha > 0 \Rightarrow D_{\text{R}, \alpha}(f_1, f_2) \geq 0, D_{\text{T}, \alpha}(f_1, f_2) \geq 0$.

– $D_{\text{SM}, \alpha, \beta}(f_1, f_2)$: $D_{\text{SM}, \alpha, \beta}(f_1, f_2) = 0$, if and only if $f_1 = f_2$.

$$D_{\text{SM}, \alpha, \beta}(f_1, f_2) = \frac{1}{\beta - 1} \left([D_{\text{temp1}}(\alpha)]^{\frac{1-\beta}{1-\alpha}} - 1 \right), \quad D_{\text{temp1}}(\alpha) = \int_{\mathbb{R}^d} [f_1(\mathbf{u})]^\alpha [f_2(\mathbf{u})]^{1-\alpha} d\mathbf{u}. \quad (44)$$

$D_{\text{temp1}}(\alpha)$ is the α -divergence, or for $\alpha = \frac{1}{2}$ the Bhattacharyya coefficient (also called Bhattacharyya kernel, or Hellinger affinity; a specific case of probability product kernels, see (59)): $BC = \int_{\mathbb{R}^d} \sqrt{f_1(\mathbf{u})} \sqrt{f_2(\mathbf{u})} d\mathbf{u} \in [0, 1]$.

D_{temp1} is a specific case of $D_{\text{temp2}}(a, b) = \int_{\mathbb{R}^d} [f_1(\mathbf{u})]^a [f_2(\mathbf{u})]^b f_1(\mathbf{u}) d\mathbf{u}$, ($a, b \in \mathbb{R}$); several divergences can be expressed by these quantities.

– D_{H}^2 : is an f-divergence [(31)] with $f(t) = \frac{1}{2}(\sqrt{t} - 1)^2$.

- $D_{L_2}(f_1, f_2)$ is non-negative, and is zero iff $f_1 = f_2$.
- D_f : $D_f(f_1, f_2) \geq 0$ with equality iff $f_1 = f_2$.
- D_{MMD} :
 - * MMD is a specific case of integral probability metrics: $D_{\text{MMD}}(f_1, f_2) = \sup_{g \in \mathcal{B}} (\mathbb{E}_{\mathbf{y}^1 \sim f_1} [g(\mathbf{y}^1)] - \mathbb{E}_{\mathbf{y}^2 \sim f_2} [g(\mathbf{y}^2)])$ with $\mathcal{B} := \{g : \|g\|_{\mathcal{H}(k)} \leq 1\}$ being the unit ball in $\mathcal{H}(k)$.
 - * MMD can be defined on topological spaces.
 - * It also acts as a ‘divergence’ on the joint and the product of the marginals in HSIC (similarly to the well-known Kullback-Leibler divergence and its extensions, see Eqs. (7) - (11)):

$$I_{\text{HSIC}}(\mathbf{y}^1, \mathbf{y}^2) = D_{\text{MMD}}(f, f_1 f_2), \quad [\mathbf{y}^1; \mathbf{y}^2] \sim f.$$

- D_{EnDist} :
 - * For $\rho(u, v) = |u - v|$, D_{EnDist} is twice the Cramer-Von Mises distance.
 - * The construction holds for (\mathcal{Z}, ρ) semimetric spaces of negative type.
 - * $D_{\text{EnDist}}(f_1, f_2) \geq 0$. $D_{\text{EnDist}}(f_1, f_2) = 0 \Leftrightarrow f_1 = f_2$ for strictly negative spaces (such as \mathbb{R}^d).
- $D_{\text{SB}, \alpha}$: For $\alpha = 2$, $[D_L(f_1, f_2)]^2 = D_{\text{NB}, 2}(f_1, f_2) = D_{\text{SB}, 2}(f_1, f_2)$.
- D_K, D_L : They are
 - * non-negative, and are zero iff $f_1 = f_2$.
 - * closely related to the Jensen-Shannon divergence in case of uniform weighting, see Eq. (45).
- D_{JS}^π :
 - * $0 \leq D_{\text{JS}}^\pi(f_1, f_2) \leq \log(2)$, $D_{\text{JS}}^\pi(f_1, f_2) = 0 \Leftrightarrow f_1 = f_2$.
 - * Specifically, for $\pi_1 = \pi_2 = \frac{1}{2}$ we obtain

$$\begin{aligned} D_{\text{JS}}(f_1, f_2) &= D_{\text{JS}}^{(\frac{1}{2}, \frac{1}{2})}(f_1, f_2) \\ &= H\left(\frac{\mathbf{y}^1 + \mathbf{y}^2}{2}\right) - \frac{H(\mathbf{y}^1) + H(\mathbf{y}^2)}{2} = \frac{1}{2} \left[D\left(f_1, \frac{f_1 + f_2}{2}\right) + D\left(f_2, \frac{f_1 + f_2}{2}\right) \right]. \end{aligned} \quad (45)$$

- $D_{\text{JT}, \alpha}$: $\lim_{\alpha \rightarrow 1} D_{\text{JT}, \alpha}(f_1, f_2) = D_{\text{JS}}(f_1, f_2)$.

3.1.4 Association Measure

- Notation: $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M] \in \mathbb{R}^d$ ($\mathbf{y}^m \in \mathbb{R}^{d_m}$, $d = \sum_{m=1}^M d_m$). F : cdf (cumulative density function) of $\mathbf{y} = [y^1; \dots; y^d]$, F_i : cdf of y^i . C : copula of \mathbf{y} , i.e. $F(\mathbf{y}) = C(F_1(y^1), \dots, F_d(y^d))$, in other words $C(\mathbf{u}) = \mathbb{P}(\mathbf{U} \leq \mathbf{u})$ where $\mathbf{U} = [F_1(y^1); \dots; F_d(y^d)] \in [0, 1]^d$. Bar stands for the survival function of its argument (it is *not* a copula in general): $\bar{C}(\mathbf{u}) := \mathbb{P}(\mathbf{U} > \mathbf{u})$. C_{kl} : bivariate marginal copula of y_{kl} , $\Pi(u_1, \dots, u_d) = \prod_{i=1}^d u_i$ (product copula), $M(\mathbf{u}) = \min_{i=1, \dots, d} u_i$ (comonotonicity copula)¹². The name of M originates from the fact that for any C copula

$$W(\mathbf{u}) := \max(u_1 + \dots + u_d - d + 1, 0) \leq C(\mathbf{u}) \leq M(\mathbf{u}), \quad \forall \mathbf{u} \in [0, 1]^d. \quad (46)$$

The well-known Spearman’s ρ (also called Spearman’s rank correlation coefficient or grade correlation coefficient) is

$$\begin{aligned} A_\rho(y^1, y^2) &= \text{corr}(F_1(y^1), F_2(y^2)) \\ &= A_\rho(C) = \frac{\int_{[0,1]^2} u_1 u_2 dC(\mathbf{u}) - (\frac{1}{2})^2}{\frac{1}{12}} = 12 \int_{[0,1]^2} C(\mathbf{u}) d\mathbf{u} - 3 = \frac{\int_{[0,1]^2} C(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^2} \Pi(\mathbf{u}) d\mathbf{u}}{\int_{[0,1]^2} M(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^2} \Pi(\mathbf{u}) d\mathbf{u}}, \end{aligned}$$

where $\int_{[0,1]^2} M(\mathbf{u}) d\mathbf{u} = \frac{1}{3}$, $\int_{[0,1]^2} \Pi(\mathbf{u}) d\mathbf{u} = \frac{1}{4}$. A_ρ can be viewed as the normalized average difference of the copula of \mathbf{y} (C) and the independence copula (Π).

- Goal: Our aim is to estimate the association of the d_m -dimensional components (\mathbf{y}^m) of the random variable $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M] \in \mathbb{R}^d$ from which we have i.i.d. samples $\{\mathbf{y}_t\}_{t=1}^T$ ($d = \sum_{m=1}^M d_m$, $\mathbf{y}^m \in \mathbb{R}^{d_m}$).

¹²Notice that in this section ‘ M ’ stands for the comonitonicity copula with argument \mathbf{u} [see $M(\mathbf{u})$] and also for the number of subspaces in $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M]$ as a subscript.

- Definition: The Spearman's ρ multivariate-1 (A_{ρ_1}), Spearman's ρ multivariate-2 (A_{ρ_2}), Spearman's ρ multivariate-3 (A_{ρ_3} ; average of A_{ρ_1} and A_{ρ_2}), Spearman's ρ multivariate-4 (A_{ρ_4} ; average pairwise Spearman's ρ), multivariate extension of Blomqvist's β (A_β ; medial correlation coefficient), multivariate conditional version of Spearman's ρ (lower tail: $A_{\rho_{\text{lt}}}$, upper tail: $A_{\rho_{\text{ut}}}$), lower and upper tail dependencies via conditional Spearman's ρ (A_{ρ_L}, A_{ρ_U}) are defined as

$$A_{\rho_1}(y^1, \dots, y^d) = A_{\rho_1}(C) = \frac{\int_{[0,1]^d} C(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^d} \Pi(\mathbf{u}) d\mathbf{u}}{\int_{[0,1]^d} M(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^d} \Pi(\mathbf{u}) d\mathbf{u}} = h_\rho(d) \left[2^d \int_{[0,1]^d} C(\mathbf{u}) d\mathbf{u} - 1 \right], \quad (47)$$

$$h_\rho(d) = \frac{d+1}{2^d - (d+1)}, \quad (48)$$

$$A_{\rho_2}(y^1, \dots, y^d) = A_{\rho_2}(C) = \frac{\int_{[0,1]^d} \Pi(\mathbf{u}) dC(\mathbf{u}) - \int_{[0,1]^d} \Pi(\mathbf{u}) d\mathbf{u}}{\int_{[0,1]^d} M(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^d} \Pi(\mathbf{u}) d\mathbf{u}} = h_\rho(d) \left[2^d \int_{[0,1]^d} \Pi(\mathbf{u}) dC(\mathbf{u}) - 1 \right], \quad (49)$$

$$A_{\rho_3}(y^1, \dots, y^d) = A_{\rho_3}(C) = \frac{A_{\rho_1}(y^1, \dots, y^d) + A_{\rho_2}(y^1, \dots, y^d)}{2}, \quad (50)$$

$$A_{\rho_4}(y^1, \dots, y^d) = A_{\rho_4}(C) = h_\rho(2) \left[2^2 \binom{d}{2}^{-1} \sum_{k,l=1; k<l}^d \int_{[0,1]^2} C_{kl}(u,v) du dv - 1 \right] = \binom{d}{2}^{-1} \sum_{k,l=1; k<l}^d A_\rho(y^k, y^l), \quad (51)$$

$$A_\beta(y^1, \dots, y^d) = A_\beta(C) = \frac{C(\mathbf{1}/2) - \Pi(\mathbf{1}/2) + \bar{C}(\mathbf{1}/2) - \bar{\Pi}(\mathbf{1}/2)}{M(\mathbf{1}/2) - \Pi(\mathbf{1}/2) + \bar{M}(\mathbf{1}/2) - \bar{\Pi}(\mathbf{1}/2)} = h_\beta(d) [C(\mathbf{1}/2) + \bar{C}(\mathbf{1}/2) - 2^{1-d}], \quad (52)$$

$$h_\beta(d) = \frac{2^{d-1}}{2^{d-1} - 1},$$

$$A_{\rho_{\text{lt}}}(y^1, \dots, y^d) = A_{\rho_{\text{lt}}}(C) = \frac{\int_{[0,p]^d} C(\mathbf{u}) d\mathbf{u} - \int_{[0,p]^d} \Pi(\mathbf{u}) d\mathbf{u}}{\int_{[0,p]^d} M(\mathbf{u}) d\mathbf{u} - \int_{[0,p]^d} \Pi(\mathbf{u}) d\mathbf{u}} = \frac{\int_{[0,p]^d} C(\mathbf{u}) d\mathbf{u} - \left(\frac{p^2}{2}\right)^d}{\frac{p^{d+1}}{d+1} - \left(\frac{p^2}{2}\right)^d}, \quad (53)$$

$$A_{\rho_{\text{ut}}}(y^1, \dots, y^d) = A_{\rho_{\text{ut}}}(C) = \frac{\int_{[1-p,1]^d} C(\mathbf{u}) d\mathbf{u} - \int_{[1-p,1]^d} \Pi(\mathbf{u}) d\mathbf{u}}{\int_{[1-p,1]^d} M(\mathbf{u}) d\mathbf{u} - \int_{[1-p,1]^d} \Pi(\mathbf{u}) d\mathbf{u}}, \quad (54)$$

$$A_{\rho_L}(y^1, \dots, y^d) = A_{\rho_L}(C) = \lim_{p \rightarrow 0, p > 0} A_{\rho_{\text{lt}}}(C) = \lim_{p \rightarrow 0, p > 0} \frac{d+1}{p^{d+1}} \int_{[0,p]^d} C(\mathbf{u}) d\mathbf{u}, \quad (55)$$

$$A_{\rho_U}(y^1, \dots, y^d) = A_{\rho_U}(C) = \lim_{p \rightarrow 0, p > 0} A_{\rho_{\text{ut}}}(C), \quad (56)$$

where $\mathbf{1}/2 = [\frac{1}{2}; \dots; \frac{1}{2}] \in \mathbb{R}^d$.

- Note:

- $A_{\rho_1}, A_{\rho_2}, A_{\rho_3}, A_{\rho_4}$: They are generalizations of Spearman's ρ , in other words $A_\rho = A_{\rho_1} = A_{\rho_2} = A_{\rho_3}$ for $d = 2$.
- A_{ρ_3}, A_{ρ_4} : These quantities are multivariate measures of concordance (see Def. 3 in Section C).
- $A_{\rho_1}, A_{\rho_2}, A_\beta$: They satisfy all the axioms of multivariate measure of concordance except for Duality.
- $A_{\rho_{\text{lt}}}, A_{\rho_{\text{ut}}}$: They belong to the following class of association measures parameterized by a function g

$$A_{\rho_g}(y^1, \dots, y^d) = A_{\rho_g}(C) = \frac{\int_{[0,1]^d} C(\mathbf{u}) g(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^d} \Pi(\mathbf{u}) g(\mathbf{u}) d\mathbf{u}}{\int_{[0,1]^d} M(\mathbf{u}) g(\mathbf{u}) d\mathbf{u} - \int_{[0,1]^d} \Pi(\mathbf{u}) g(\mathbf{u}) d\mathbf{u}}.$$

- $A_{\rho_{\text{lt}}}$:

- * Here $g(\mathbf{u}) = \mathbb{I}_{[0,p]^d}(\mathbf{u})$, where $0 < p \leq 1$ and \mathbb{I} is the indicator function. This g choice refers to the weighting of the lower part of the copula, i.e., we measure the amount of dependence in the lower tail of the multivariate distributions. For $p = 1$, $A_{\rho_{\text{lt}}} = A_{\rho_1}$.

- * It preserves the concordance ordering [see Eq. (78)], i.e., $C_1 \prec C_2 \Rightarrow A_{\rho_{\text{lt}}}(C_1) \leq A_{\rho_{\text{lt}}}(C_2)$, for $\forall p \in (0, 1]$. Thus, from $C \prec M$ [see Eq. (46)] one obtains that $A_{\rho_{\text{lt}}} \leq 1$.

- $A_{\rho_{\text{ut}}}$: In this case $g(\mathbf{u}) = \mathbb{I}_{[1-p,1]^d}(\mathbf{u})$, where $0 < p \leq 1$; in other words the weighting is put on the upper tail.

3.1.5 Cross Quantity

- Notation: $\mathbb{R}^d \ni \mathbf{y}^1 \sim f_1, \mathbb{R}^d \ni \mathbf{y}^2 \sim f_2$.
- Goal: We want to estimate cross quantities from independent, i.i.d. samples $\{\mathbf{y}_t^1\}_{t=1}^{T_1}$ and $\{\mathbf{y}_t^2\}_{t=1}^{T_2}$ distributed according to f_1 and f_2 , respectively.
- Definition: The cross-entropy (C_{CE}) is

$$C_{CE}(f_1, f_2) = - \int_{\mathbb{R}^d} f_1(\mathbf{u}) \log [f_2(\mathbf{u})] d\mathbf{u}. \quad (57)$$

3.1.6 Kernel on Distributions

- Notation: $\mathbb{R}^d \ni \mathbf{y}^1 \sim f_1, \mathbb{R}^d \ni \mathbf{y}^2 \sim f_2$.
- Goal: Our aim is to estimate the value of a kernel $[K(f_1, f_2)]$ given independent, i.i.d. samples from \mathbf{y}^1 and \mathbf{y}^2 , $\{\mathbf{y}_t^1\}_{t=1}^{T_1}$ and $\{\mathbf{y}_t^2\}_{t=1}^{T_2}$.
- Definition: the expected kernel (K_{exp} ; also called summation kernel, mean map kernel, set kernel, multi-instance kernel, ensemble kernel; a specific convolution kernel), probability product kernel ($K_{PP,\rho}$), Jensen-Shannon kernel (K_{JS}), Jensen-Tsallis kernel ($K_{JT,\alpha}$), exponentiated Jensen-Shannon kernel ($K_{EJS,u}$), exponentiated Jensen-Rényi kernels ($K_{EJR1,u,\alpha}$, $K_{EJR2,u,\alpha}$), exponentiated Jensen-Tsallis kernels ($K_{EJT1,u,\alpha}$, $K_{EJT2,u,\alpha}$) are defined as

$$K_{\text{exp}}(f_1, f_2) = \langle \mu_1, \mu_2 \rangle_{\mathcal{H}(k)} = \mathbb{E}_{\mathbf{y}^1, \mathbf{y}^2} [k(\mathbf{y}^1, \mathbf{y}^2)], \quad \mu_i = \mathbb{E}_{\mathbf{y}^i \sim f_i} [k(\cdot, \mathbf{y}^i)], \quad (58)$$

$$K_{PP,\rho}(f_1, f_2) = \int_{\mathbb{R}^d} [f_1(\mathbf{u})]^\rho [f_2(\mathbf{u})]^\rho d\mathbf{u}, \quad (\rho > 0), \quad (59)$$

$$K_{JS}(f_1, f_2) = \log(2) - D_{JS}(f_1, f_2), \quad (60)$$

$$K_{JT,\alpha}(f_1, f_2) = \log_\alpha(2) - T_\alpha(f_1, f_2), \quad (\alpha \in (0, 2] \setminus \{1\}), \quad \log_\alpha(x) = \frac{x^{1-\alpha} - 1}{1-\alpha}, \quad (61)$$

$$T_\alpha(f_1, f_2) = H_{T,\alpha} \left(\frac{\mathbf{y}^1 + \mathbf{y}^2}{2} \right) - \frac{H_{T,\alpha}(\mathbf{y}^1) + H_{T,\alpha}(\mathbf{y}^2)}{2^\alpha},$$

$$K_{EJS,u}(f_1, f_2) = e^{-uD_{JS}(f_1, f_2)}, \quad (u > 0), \quad (62)$$

$$K_{EJR1,u,\alpha}(f_1, f_2) = e^{-uH_{R,\alpha} \left(\frac{\mathbf{y}^1 + \mathbf{y}^2}{2} \right)}, \quad (u > 0, \alpha \in (0, 1)), \quad (63)$$

$$K_{EJR2,u,\alpha}(f_1, f_2) = e^{-uD_{JR,\alpha}(f_1, f_2)}, \quad (u > 0, \alpha \in (0, 1)) \quad D_{JR,\alpha}(f_1, f_2) = D_{JR,\alpha}^{\left(\frac{1}{2}, \frac{1}{2}\right)}(f_1, f_2), \quad (64)$$

$$K_{EJT1,u,\alpha}(f_1, f_2) = e^{-uH_{T,\alpha} \left(\frac{\mathbf{y}^1 + \mathbf{y}^2}{2} \right)}, \quad (u > 0, \alpha \in (0, 2] \setminus \{1\}), \quad (65)$$

$$K_{EJT2,u,\alpha}(f_1, f_2) = e^{-uD_{JT,\alpha}(f_1, f_2)}, \quad (u > 0, \alpha \in (0, 2] \setminus \{1\}). \quad (66)$$

- Note:

– K_{exp} : it

* generates MMD, $[D_{\text{MMD}}(f_1, f_2)]^2 = K_{\text{exp}}(f_1, f_1) - 2K_{\text{exp}}(f_1, f_2) + K_{\text{exp}}(f_2, f_2)$.

* can be defined slightly more generally, on topological spaces.

– $K_{PP,\rho}$: For $\rho = \frac{1}{2}$ we get back the Bhattacharyya kernel (K_B ; also known as Bhattacharyya coefficient, or Hellinger affinity) K_B is intimately related to, induces the Hellinger distance [see Eq. (29)]:

$$[D_H(f_1, f_2)]^2 = \frac{1}{2} [K_B(f_1, f_1) - 2K_B(f_1, f_2) + K_B(f_2, f_2)] = \frac{1}{2} [2 - 2K_B(f_1, f_2)] = 1 - K_B(f_1, f_2). \quad (67)$$

– $K_{JT,\alpha}$: $\lim_{\alpha \rightarrow 1} K_{JT,\alpha}(f_1, f_2) = K_{JS}(f_1, f_2)$.

– $K_{EJR2,u,\alpha}$: $\lim_{\alpha \rightarrow 1} K_{EJR2,u,\alpha}(f_1, f_2) = K_{EJS,u}(f_1, f_2)$.

– $K_{EJT2,u,\alpha}$: $\lim_{\alpha \rightarrow 1} K_{EJT2,u,\alpha}(f_1, f_2) = K_{EJS,u}(f_1, f_2)$.

Estimated quantity	Principle	d	Cost name
Shannon entropy (H)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BHShannon_KnnK
Shannon entropy (H)	approximate slope of the inverse distribution function	$= 1$	BHShannon_SpacingV
Shannon entropy (H)	maximum entropy distribution, function set1, plug-in	$= 1$	BHShannon_MaxEnt1
Shannon entropy (H)	maximum entropy distribution, function set2, plug-in	$= 1$	BHShannon_MaxEnt2
Rényi entropy ($H_{R,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BHRenyi_KnnK
Rényi entropy ($H_{R,\alpha}$)	generalized nearest neighbors ($S \subseteq \{1, \dots, k\}$)	≥ 1	BHRenyi_KnnS
Tsallis entropy ($H_{T,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BHTsallis_KnnK
Sharma-Mittal entropy ($H_{SM,\alpha,\beta}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BHSharmaMittal_KnnK
Φ -entropy ($H_{\Phi,w}$)	sample spacing	$= 1$	BHPhi_Spacing
Shannon entropy (H)	-KL divergence from the normal distribution: (70)	≥ 1	MHShannon_DKLN
Shannon entropy (H)	-KL divergence from the uniform distribution: (71)	≥ 1	MHShannon_DKLU
Tsallis entropy ($H_{T,\alpha}$)	function of the Rényi entropy: (3)	≥ 1	MHTsallis_HR

Table 1: Entropy estimators. Third column: dimension (d) constraint. Top: base methods, bottom: meta estimators.

3.2 Conditional Quantities

The toolbox supports the estimation of conditional quantities defined in Section 3.2.1 (entropy) and Section 3.2.2 (mutual information).

3.2.1 Entropy

- Notation: $\mathbf{y} = [\mathbf{y}^1; \mathbf{y}^2]$, $\mathbf{y}^m \in \mathbb{R}^{d_m}$.
- Goal: Assume we have $\{(\mathbf{y}_t^1, \mathbf{y}_t^2)\}_{t=1}^T$ samples; we want to estimate the conditional entropy of \mathbf{y}^1 given \mathbf{y}^2 .
- Definition: The conditional Shannon entropy $[H(\cdot|\cdot)]$ is defined as

$$H(\mathbf{y}^1|\mathbf{y}^2) = \mathbb{E}_{\mathbf{y}^2} [H(\mathbf{y}^1|\mathbf{y}^2)] = H([\mathbf{y}^1; \mathbf{y}^2]) - H(\mathbf{y}^2). \quad (68)$$

3.2.2 Mutual Information

- Notation: $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M; \mathbf{y}^{M+1}] \in \mathbb{R}^d$, where $\mathbf{y}^m \in \mathbb{R}^{d_m}$ and $d = \sum_{m=1}^{M+1} d_m$.
- Goal: Assume we have access to the samples $\{(\mathbf{y}_t^1, \dots, \mathbf{y}_t^M, \mathbf{y}_t^{M+1})\}_{t=1}^T$, our aim is to estimate the mutual information of $\mathbf{y}^1, \dots, \mathbf{y}^M$ given \mathbf{y}^{M+1} .
- Definition: The conditional Shannon mutual information $[I(\cdot|\cdot)]$ is

$$\begin{aligned} I(\mathbf{y}^1, \dots, \mathbf{y}^M | \mathbf{y}^{M+1}) &= \mathbb{E}_{\mathbf{y}^{M+1}} [I(\mathbf{y}^1, \dots, \mathbf{y}^M | \mathbf{y}^{M+1})] \\ &= -H([\mathbf{y}^1; \dots; \mathbf{y}^{M+1}]) + \sum_{m=1}^M H([\mathbf{y}^m; \mathbf{y}^{M+1}]) - (M-1)H(\mathbf{y}^{M+1}). \end{aligned} \quad (69)$$

3.3 Estimators

The available estimators in ITE are listed in

1. **unconditional quantities:** Table 1 (entropy), Table 2 (mutual information), Table 3 (divergence), Table 4 (association measures), Table 5 (cross quantity), Table 6 (kernel on distributions).
2. **conditional quantities:** Table 7 (entropy), Table 8 (mutual information).

Demos of the estimators are enlisted in Table 9.

- Certain association measures and mutual information estimators require one-dimensional subspaces ($\forall d_m = 1$). For these estimators the toolbox provides a simplified calling syntax (without specifying $\{d_m\}_{m=1}^M$):

Estimated quantity	Principle	d_m	M	Cost name
kernel canonical correlation (I_{KCCA})	sup correlation over RKHSs	≥ 1	≥ 2	BIKGV
kernel generalized variance (I_{KGV})	Gaussian mutual information of the features	≥ 1	≥ 2	BIKCCA
Hoeffding's Φ (I_Φ), multivariate	L^2 distance of the joint- and the product copula	$= 1$	≥ 2	BIHoeffding
Hilbert-Schmidt indep. criterion (I_{HSIC})	HS norm of the cross-covariance operator	≥ 1	≥ 2	BIHSIC_IChol
distance covariance (I_{dCov})	pairwise distances	≥ 1	$= 2$	BIDistCov
distance correlation (I_{dCor})	pairwise distances	≥ 1	$= 2$	BIDistCorr
Lancaster 3-variable interaction (I_{3-Lanc})	embedding of the Lancaster interaction measure	≥ 1	$= 3$	BI3WayJoint
3-variable joint independence ($I_{3-joint}$)	embedding of the 'joint - product of marginals'	≥ 1	$= 3$	BI3WayLancaster
(Shannon) mutual information (I)	KL-divergence of joint & product of marginals: (7)	≥ 1	≥ 2	MIShannon_DKL
(Shannon) mutual information (I)	entropy sum of components minus joint entropy: (72)	≥ 1	≥ 2	MIShannon_HS
Rényi mutual information ($I_{R,\alpha}$)	Rényi divergence of joint & product of marginals: (8)	≥ 1	≥ 2	MIRenyi_DR
Rényi mutual information ($I_{R,\alpha}$)	minus the Rényi entropy of the joint copula: (73)	$= 1$	≥ 2	MIRenyi_HR
Tsallis mutual information ($I_{T,\alpha}$)	Tsallis divergence of joint & product of marginals: (9)	≥ 1	≥ 2	MITsallis_DT
χ^2 mutual information (I_{χ^2})	χ^2 divergence of joint & product of marginals: (10)	≥ 1	≥ 2	MIChi2_DChi2
L_2 mutual information (I_{L_2})	L_2 -divergence of joint & product of marginals: (11)	≥ 1	≥ 2	MIL2_DL2
copula-based kernel dependency (I_c)	MMD div. of the joint copula & uniform distr.: (12)	$= 1$	≥ 2	MIMMD_CopulaDMMD
distance covariance (I_{dCov})	pairwise distances, equivalence to HSIC: (23)	≥ 1	$= 2$	MIDistCov_HSIC

Table 2: Mutual information estimators. Third column: dimension constraint (d_m ; $\mathbf{y}^m \in \mathbb{R}^{d_m}$). Fourth column: constraint for the number of components (M ; $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M]$). Top: base methods, bottom: meta estimators.

Estimated quantity	Principle	d	Cost name
Kullback-Leibler divergence (D)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDKL_KnnK
Kullback-Leibler divergence (D)	k-nearest neighbors ($S_i = \{k_i(T_i)\}$)	≥ 1	BDKL_KnnKiTi
Rényi divergence ($D_{R,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDRenyi_KnnK
Tsallis divergence ($D_{T,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDTsallis_KnnK
Sharma-Mittal divergence ($D_{SM,\alpha,\beta}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDSharmaMittal_KnnK
Pearson χ^2 divergence (D_{χ^2})	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDChi2_KnnK
Hellinger distance (D_H)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDEllinger_KnnK
L_2 divergence (D_{L_2})	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDL2_KnnK
maximum mean discrepancy (D_{MMD})	U-statistic, unbiased	≥ 1	BDMMD_UStat
maximum mean discrepancy (D_{MMD})	V-statistic, biased	≥ 1	BDMMD_VStat
maximum mean discrepancy (D_{MMD})	U-statistic, incomplete Cholesky decomposition	≥ 1	BDMMD_UStat_IChol
maximum mean discrepancy (D_{MMD})	V-statistic, incomplete Cholesky decomposition	≥ 1	BDMMD_VStat_IChol
maximum mean discrepancy (D_{MMD})	online	≥ 1	BDMMD_Online
energy distance (D_{EnDist})	pairwise distances	≥ 1	BDEnergyDist
Bhattacharyya distance (D_B)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDBhattacharyya_KnnK
Bregman distance ($D_{NB,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDBregman_KnnK
symmetric Bregman distance ($D_{SB,\alpha}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BDSymBregman_KnnK
Kullback-Leibler divergence (D)	difference of cross-entropy and entropy: (74)	≥ 1	MDKL_HSCE
f-divergence (D_f)	second-order Taylor expansion, χ^2 divergence: (75)	≥ 1	MDf_DChi2
maximum mean discrepancy (D_{MMD})	block-average of U-statistic based MMDs	≥ 1	MDBlockMMD
energy distance (D_{EnDist})	pairwise distances, equivalence to MMD: (76)	≥ 1	MDEnergyDist_DMMD
symmetric Bregman distance ($D_{SB,\alpha}$)	symmetrised Bregman distance: (36)	≥ 1	MDSymBregman_DB
J-distance (D_J)	symmetrised Kullback-Leibler divergence: (38)	≥ 1	MDJDist_DKL
K divergence (D_K)	smoothed Kullback-Leibler divergence: (39)	≥ 1	MDK_DKL
L divergence (D_L)	symmetrised K divergence: (40)	≥ 1	MDL_DKL
Jensen-Shannon divergence (D_{JS}^π)	smoothed (π), defined via the Shannon entropy: (41)	≥ 1	MDJS_HS
Jensen-Rényi divergence ($D_{JR,\alpha}^\pi$)	smoothed (π), defined via the Rényi entropy: (42)	≥ 1	MDJR_HR
Jensen-Tsallis divergence ($D_{JT,\alpha}$)	smoothed, defined via the Tsallis entropy: (43)	≥ 1	MDJT_HT

Table 3: Divergence estimators. Third column: dimension (d) constraint. Top: base methods, bottom: meta estimators.

Estimated quantity	Principle	d_m	M	Cost name
Spearman's ρ : multivariate1 (A_{ρ_1})	empirical copula, explicit formula	= 1	≥ 2	BASpearman1
Spearman's ρ : multivariate2 (A_{ρ_2})	empirical copula, explicit formula	= 1	≥ 2	BASpearman2
Spearman's ρ : multivariate3 (A_{ρ_3})	average of ρ_1 and ρ_2	= 1	≥ 2	BASpearman3
Spearman's ρ : multivariate4 (A_{ρ_4})	average pairwise Spearman's ρ	= 1	≥ 2	BASpearman4
Blomqvist's β (A_β)	empirical copula, explicit formula	= 1	≥ 2	BABlomqvist
conditional Spearman's ρ , lower tail ($A_{\rho_{lt}}$)	empirical copula, explicit formula	= 1	≥ 2	BASpearmanCondLT
conditional Spearman's ρ , upper tail ($A_{\rho_{ut}}$)	empirical copula, explicit formula	= 1	≥ 2	BASpearmanCondUT
lower tail dep. via conditional Spearman's ρ (A_{ρ_L})	limit of $A_{\rho_{lt}}$: (55)	= 1	≥ 2	MASpearmanLT
upper tail dep. via conditional Spearman's ρ (A_{ρ_U})	limit of $A_{\rho_{ut}}$: (56)	= 1	≥ 2	MASpearmanUT

Table 4: Association measure estimators. Third column: dimension constraint (d_m ; $\mathbf{y}^m \in \mathbb{R}^{d_m}$). Fourth column: constraint for the number of components (M ; $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M]$). Top: base methods, bottom: meta estimators.

Estimated quantity	Principle	d	Cost name
cross-entropy (C_{CE})	k-nearest neighbors ($S = \{k\}$)	≥ 1	BCCE_KnnK

Table 5: Cross quantity estimators. Third column: dimension (d) constraint.

Estimated quantity	Principle	d	Cost name
expected kernel (K_{exp})	mean of pairwise kernel values	≥ 1	BKExpected
probability product kernel ($K_{PP,\rho}$)	k-nearest neighbors ($S = \{k\}$)	≥ 1	BKProbProd_KnnK
Jensen-Shannon kernel (K_{JS})	function of the Jensen-Shannon divergence: (60)	≥ 1	MKJS_DJS
Jensen-Tsallis kernel ($K_{JT,\alpha}$)	function of the Tsallis entropy: (61)	≥ 1	MKJT_HT
exponentiated Jensen-Shannon kernel ($K_{EJS,u}$)	function of the Jensen-Shannon divergence: (62)	≥ 1	MKExpJS_DJS
exponentiated Jensen-Rényi kernel-1 ($K_{EJR1,u,\alpha}$)	function of the Rényi entropy: (63)	≥ 1	MKExpJR1_HR
exponentiated Jensen-Rényi kernel-2 ($K_{EJR2,u,\alpha}$)	function of the Jensen-Rényi divergence: (64)	≥ 1	MKExpJR2_DJR
exponentiated Jensen-Tsallis kernel-1 ($K_{EJT1,u,\alpha}$)	function of the Tsallis entropy: (65)	≥ 1	MKExpJT1_HT
exponentiated Jensen-Tsallis kernel-2 ($K_{EJT2,u,\alpha}$)	function of the Jensen-Tsallis divergence: (66)	≥ 1	MKExpJT2_DJT

Table 6: Estimators of kernels on distributions. Third column: dimension (d) constraint. Top: base methods, bottom: meta estimators.

Estimated quantity	Principle	d_m	Cost name
conditional Shannon entropy [$H(\cdot \cdot)$]	reduction to Shannon entropy	≥ 1	BcondHShannon_HShannon

Table 7: Conditional entropy estimators. Third column: dimension (d_m) constraint.

Estimated quantity	Principle	d_m	M	Cost name
conditional Shannon mutual information [$I(\cdot \cdot)$]	reduction to Shannon entropy	≥ 1	≥ 2	BcondIShannon_HShannon

Table 8: Conditional mutual information estimators. Third column: dimension constraint (d_m ; $\mathbf{y}^m \in \mathbb{R}^{d_m}$). Fourth column: constraint for the number of components (M ; $\mathbf{y} = [\mathbf{y}^1; \dots; \mathbf{y}^M; \mathbf{y}^{M+1}]$).

```
>>> a_or_i      = co. estimation(y, ds) # long (traditional) syntax
>>> a_or_i_v2  = co. estimation(y)    # short syntax; gives the same result as the previous line
```

- We use the `TypeXName` naming convention.
 - For example, `BHShannon_KnnK` = `B` + `H` + `Shannon_KnnK`. `B` means base estimator, `H` is for entropy, `Shannon_KnnK` is about the estimated quantity/technique (Shannon entropy, with `k`-nearest neighbors).
 - Generally, `Type` \in `{B,M}`, `X` \in `{H, I, D, A, C, K, condH, condI}`, where `B`-base, `M`-meta; `H`-entropy, `I`-mutual information, `D`-divergence, `A`-association measure, `C`-cross quantity, `K`-kernel on distributions, `condH`-conditional entropy, `condI`-conditional mutual information.
- Meta estimators: The rules not yet covered for the meta estimators are as follows.
 - Notation: $\mathbf{y} \sim f$; f_U : uniform density on $[0, 1]^d$; $N(\mathbf{m}, \Sigma)$: normal distribution with mean \mathbf{m} and covariance matrix Σ . $cov(\mathbf{y})$ denotes the covariance of \mathbf{y} .
 - Rules:

$$H(\mathbf{y}) = H(\mathbf{y}_G) - D(f, f_G) \qquad \mathbf{y}_G \sim f_G = N(\mathbb{E}(\mathbf{y}), cov(\mathbf{y})), \quad (70)$$

$$H(\mathbf{y}) = -D(f, f_U) \qquad \mathbf{y} \in [0, 1]^d \text{ (if } \mathbf{y} \in [\mathbf{a}, \mathbf{b}] = \times_{i=1}^d [a_i, b_i], \quad (71)$$

it is linearly transformed to $[0, 1]^d$),

$$I(\mathbf{y}^1, \dots, \mathbf{y}^M) = \sum_{m=1}^M H(\mathbf{y}^m) - H([\mathbf{y}^1; \dots; \mathbf{y}^M]), \quad (72)$$

$$I_{R,\alpha}(y^1, \dots, y^M) = -H_{R,\alpha}(\mathbf{z}), \qquad \mathbf{z} = [F_1(y^1); \dots; F_M(y^M)] \in \mathbb{R}^M, \quad (73)$$

$$D(f_1, f_2) = C_{CE}(f_1, f_2) - H(f_1), \quad (74)$$

$$D_f(f_1, f_2) \approx \frac{f''(1)}{2} D_{\chi^2}(f_1, f_2), \quad (75)$$

$$D_{EnDist}(f_1, f_2) = 2 [D_{MMD}(f_1, f_2)]^2, \qquad \rho(\mathbf{u}, \mathbf{v}) = k(\mathbf{u}, \mathbf{u}) + k(\mathbf{v}, \mathbf{v}) - 2k(\mathbf{u}, \mathbf{v}), \quad (76)$$

where D_{EnDist} is determined by ρ , D_{MMD} by k .

- Base estimators: Their equations can be looked up by using Section B.

Estimated quantity	ite/demos/
Shannon entropy (H)	demo_h_shannon.py
Rényi entropy ($H_{R,\alpha}$)	demo_h_renyi.py
Tsallis entropy ($H_{T,\alpha}$)	demo_h_tsallis.py
Sharma-Mittal entropy ($H_{SM,\alpha,\beta}$)	demo_h_sharma_mittal.py
Φ -entropy ($H_{\Phi,w}$)	demo_h_phi.py
(Shannon) mutual information (I)	demo_i_shannon.py
Rényi mutual information ($I_{R,\alpha}$)	demo_i_renyi.py
Kullback-Leibler divergence (D)	demo_d_kullback_leibler.py
Rényi divergence ($D_{R,\alpha}$)	demo_d_renyi.py
Tsallis divergence ($D_{T,\alpha}$)	demo_d_tsallis.py
Sharma-Mittal divergence ($D_{SM,\alpha,\beta}$)	demo_d_sharma_mittal.py
Pearson χ^2 divergence (D_{χ^2})	demo_d_chi_square.py
Hellinger distance (D_H)	demo_d_hellinger.py
L_2 divergence (D_{L_2})	demo_d_l2.py
Maximum mean discrepancy (D_{MMD})	demo_d_mmd.py
Bregman distance ($D_{NB,\alpha}$)	demo_d_bregman.py
Jensen-Rényi divergence ($D_{JR,\alpha}^\pi$)	demo_d_jensen_renyi.py
cross-entropy (C_{CE})	demo_c_cross_entropy.py
expected kernel (K_{exp})	demo_k_expected.py
probability product kernel ($K_{PP,\rho}$)	demo_k_prob_product.py
exponentiated Jensen-Rényi kernel-1 ($K_{EJR1,u,\alpha}$)	demo_k_ejr1.py
exponentiated Jensen-Rényi kernel-2 ($K_{EJR2,u,\alpha}$)	demo_k_ejr2.py
exponentiated Jensen-Tsallis kernel-1 ($K_{EJT1,u,\alpha}$)	demo_k_ejt1.py
exponentiated Jensen-Tsallis kernel-2 ($K_{EJT2,u,\alpha}$)	demo_k_ejt2.py
conditional Shannon entropy [$H(\cdot \cdot)$]	demo_h_shannon_cond.py
conditional Shannon mutual information [$I(\cdot \cdot)$]	demo_i_shannon_cond.py
approximation quality of incomplete Cholesky decomposition	demo_incomplete_cholesky.py
independence of y^m -s $\stackrel{?}{\Rightarrow} A(y^1, \dots, y^M) = 0$	demo_a_independence.py
independence of \mathbf{y}^m -s $\stackrel{?}{\Rightarrow} I(\mathbf{y}^1, \dots, \mathbf{y}^M) = 0$	demo_i_independence.py
$f_1 = f_2 \stackrel{?}{\Rightarrow} D(f_1, f_2) = 0$	demo_d_equality.py
$f_1, \dots, f_M \stackrel{?}{\Rightarrow} G = [G_{ij}] = [K(f_i, f_j)]_{i,j=1}^M$: positive semi-definite	demo_k_positive_semidefinite.py

Table 9: Top-middle: demos for analytical formula vs. estimated value (**analytical_values** subfolder); unconditional quantities (top), conditional ones (middle). 1st column: estimated quantity. 2nd column: .py. Bottom: independence demos, equality test, quality demo of incomplete Cholesky decomposition, positive semi-definiteness test for kernels on distributions (**other** subfolder). 1st column: task. 2nd column: .py.

A For Developers

Section [A.1](#) is about the directory structure of the package. Section [A.2](#) focuses on doctests. Adding new estimators is the topic of Section [A.3](#). Passing parameters in certain meta-estimators is detailed in Section [A.4](#).

A.1 Directory Structure

The `ite` package is organized as follows:

- `demos`: demos of the estimators.
 - `analytical_values`: analytical expressions (for a few distributions) vs. estimated values.
 - `other`: incomplete Cholesky decomposition, positive semi-definiteness of the Gram matrix defined by a kernel on distributions.
- `doc`: link to this manual.
- `ite`: contains the estimators themselves.
 - Directory `cost`:
 - * In case of the
 1. unconditional quantities: the base estimators can be found in `base_a.py` (association measures), `base_c.py` (cross quantities), `base_d.py` (divergence measures), `base_h.py` (entropy), `base_i.py` (mutual information), `base_k.py` (kernels on distributions). The meta ones are in `meta_a.py`, `meta_c.py`, `meta_d.py`, `meta_h.py`, `meta_i.py`, `meta_k.py`.
 2. conditional ones: the meta estimators are in `meta_h_cond.py` (entropy), `meta_i_cond.py` (mutual information).
 - * `x_initialization.py`, `x_verification.py`: classes to code up new estimators rapidly (initialization and verification routines).
 - * `x_factory.py`: general module to invoke estimators.
 - * `x_analytical_values.py`: analytical values for a few information theoretical quantities.
 - * `x_kernel.py`: Kernel class.
 - * `x_python_to_matlab.py`: Python ITE \leftrightarrow Matlab ITE transitions (see Section [B](#)).
 - * `__init__.py`: it makes the estimators available upon 'import ite'.
 - `__init__.py`: it loads the cost module upon 'import ite'.
 - `shared.py`: code shared by the estimators.

A.2 Running Doctests

Assumption: you have Nose installed.⁶ Change to the main `ite` folder (containing the `.txt`-s, `doc`, `ite`, ...), and issue the command

```
> nosetests --with-doctest -w ite      # run only doctests of ite/ite; 'nose' provides 'nosetests'
> nosetests --with-doctest base_a.py  # after cd-ing to ite/ite/cost, run the doctests of a single
                                     # file (base_a.py)
> nosetests --with-doctest base_a.py:BASpearman1 # doctest of a specific class/function in base_a.py
> nosetests --with-doctest base_a      # the .py extension can be discarded
> nosetests --with-doctest base_a:BASpearman1 # -||-
```

A.3 Adding New Estimators

Upon creating a new estimator (H/I/D/A/C/K/condH/condI):

1. Recall the `TypeXName` naming convention (see Section 3.3).
2. The classes in Table 10 and Table 11 can be used for initialization and verification of the estimators.

Notes:

- 'InitX' is the default base class providing printing functionality (see below), and sets `mult`. Currently, multiplicative constants are considered to be relevant; in other words the default value of 'mult' is 'True'.¹³

```
>>> import ite                                # load the ite package
>>> co = ite.cost.BHShannon_KnnK()           # initialize an entropy estimator
>>> print(co)                                 # print it(s parameters)
```

This printing capability is what we used in Example 1.

- 'VerCompSubspaceDims' is used for A/I estimators: an exception is raised if the subspace dimensions ($\{d_m\}_{m=1}^M$) and the dimension of the samples ($dim(\mathbf{y}_t)$) are not compatible.
- 'VerEqualDSubspaces' guarantees in C/D/K estimators that an exception occurs if the dimensions of the samples from \mathbf{y}^1 and \mathbf{y}^2 are different.
- 'InitBagGram' is the base class for *kernels on distributions* (empirically on bags of points), giving Gram matrix computation capability.
- Notice the two *different* meanings of 'kernel' in Table 10:
 - By 'InitKernel' one can build a kernel-based information theoretical estimator; it does not have to be a kernel on distributions. The currently implemented kernels are

$$\begin{aligned}
 k_G(a, b) &= e^{-\frac{\|a-b\|_2^2}{2\theta^2}}, & k_e(a, b) &= e^{-\frac{\|a-b\|_2}{2\theta^2}}, & k_C(a, b) &= \frac{1}{1 + \frac{\|a-b\|_2^2}{\theta^2}}, \\
 k_t(a, b) &= \frac{1}{1 + \|a-b\|_2^\theta}, & k_p(a, b) &= (\langle a, b \rangle + \theta)^p, & k_r(a, b) &= 1 - \frac{\|a-b\|_2^2}{\|a-b\|_2^2 + \theta}, \\
 k_i(a, b) &= \frac{1}{\sqrt{\|a-b\|_2^2 + \theta^2}}, & k_{M, \frac{3}{2}}(a, b) &= \left(1 + \frac{\sqrt{3}\|a-b\|_2}{\theta}\right) e^{-\frac{\sqrt{3}\|a-b\|_2}{\theta}}, \\
 k_{M, \frac{5}{2}}(a, b) &= \left(1 + \frac{\sqrt{5}\|a-b\|_2}{\theta} + \frac{5\|a-b\|_2^2}{3\theta^2}\right) e^{-\frac{\sqrt{5}\|a-b\|_2}{\theta}}.
 \end{aligned}$$

- The two meanings/capabilities can be used *independently*. For example, `MKJS_DJS` uses 'InitBagGram' but not 'InitKernel'; `BKExpected` relies on 'InitBagGram' and 'InitKernel'; `BDMMD_UStat` comes from 'InitKernel' but not from 'InitBagGram'.

3. A simplified calling syntax is provided for an A/I estimator with ' $\forall d_m = 1$ ' constraint (see Section 3.3).
4. Forcing 'mult=True' in the children [see e.g., (60), (62), (74)] and passing parameters to them (such as in Table 12) are implemented where it is necessary.
5. If the name of the estimator is added to `ite/ite/cost/_init_.py`, it is loaded automatically upon 'import ite'.

A.4 Parameter Passing for (Certain) Meta Estimators

Certain meta estimators set others' parameters during the estimation; see Table 12 for a summary.

¹³Occasionally significant computation can be saved if these multiplicative factors do not matter.

Class(Parent)	Feature
InitX(object)	string representation, initialization: mult
InitKnnK(InitX)	string representation, initialization: mult, kNN ($S = \{k\}$)
InitKnnKiTi(InitX)	string representation, initialization: mult, k-NN ($S = \{k_i(T_i)\}$)
InitAlpha(InitX)	string representation, initialization: mult, $\alpha \neq 1$
InitUAlpha(InitAlpha)	string representation, initialization: mult, $\alpha \neq 1, u > 0$
InitKnnKAlpha(InitAlpha)	string representation, initialization: mult, kNN ($S = \{k\}$), $\alpha \neq 1$
InitKnnKAlphaBeta(InitKnnKAlpha)	string representation, initialization: mult, kNN ($S = \{k\}$), $\alpha \neq 1, \beta \neq 1$
InitKnnSAlpha(InitAlpha)	string representation, initialization: mult, generalized kNN ($S \subseteq \{1, \dots, k\}$), $\alpha \neq 1$
InitKernel(InitX)	string representation, initialization: mult, kernel
InitEtaKernel(InitKernel)	string representation, initialization: mult, kernel, $\eta > 0$ (incomplete Cholesky decomposition)
InitBagGram(object)	Gram matrix computation for kernels on distributions (empirically on bags of points)

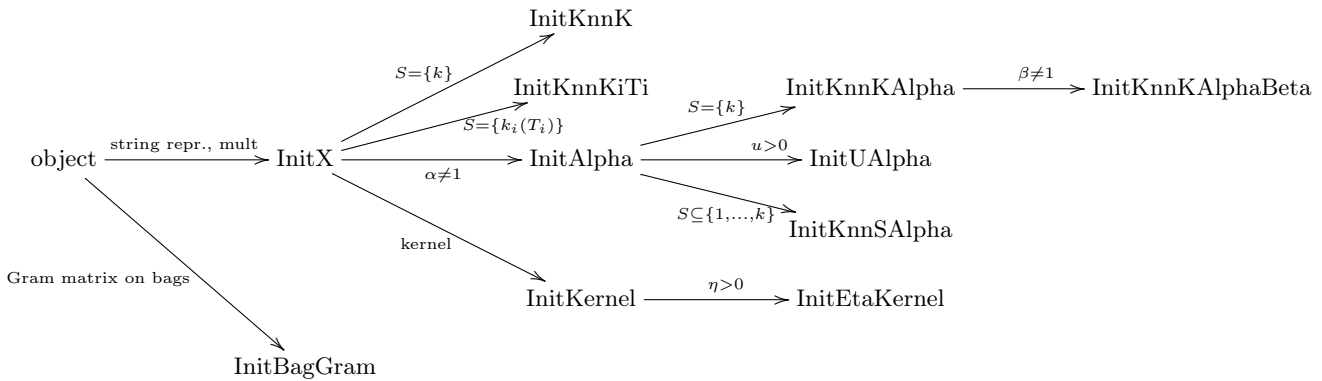


Table 10: Classes for initialization of the estimators (see `x_initialization.py`). 1st column: name of the class and its parent. 2nd column: its feature. The dependence of the classes and the additional features are schematically summarized in the diagram above.

Class	Feature
VerOneDSignal	$\mathbf{Y}_{1:T} \stackrel{?}{\in} \mathbb{R}^{T \times 1}$
VerOneDSubspaces	$d_1 = d_2 = \dots = d_M = 1?$
VerCompSubspaceDims	$[d_1; \dots; d_M], \mathbf{Y}_{1:T} \in \mathbb{R}^{T \times d} \stackrel{?}{\Rightarrow} d = \sum_{m=1}^M d_m$
VerSubspaceNumberIsK	$[d_1; \dots; d_M] \stackrel{?}{\in} \mathbb{R}^K$
VerEqualDSubspaces	$\mathbf{Y}_{1:T_1}^1 \in \mathbb{R}^{T_1 \times d_1}, \mathbf{Y}_{1:T_2}^2 \in \mathbb{R}^{T_2 \times d_2} \stackrel{?}{\Rightarrow} d_1 = d_2$
VerEqualSampleNumbers	$\mathbf{Y}_{1:T_1}^1 \in \mathbb{R}^{T_1 \times d_1}, \mathbf{Y}_{1:T_2}^2 \in \mathbb{R}^{T_2 \times d_2} \stackrel{?}{\Rightarrow} T_1 = T_2$
VerEvenSampleNumbers	$\mathbf{Y}_{1:T}^1 \in \mathbb{R}^{T \times d_1}, \mathbf{Y}_{1:T}^2 \in \mathbb{R}^{T \times d_2} \stackrel{?}{\Rightarrow} 2 T$

Table 11: Classes for verification of the estimators (see `x_verification.py`). 1st column: name of the class. 2nd column: its feature.

Inherited parameter	Eq.	Cost name
$H_{T,\alpha} \xrightarrow{\alpha} H_{R,\alpha}$	(3)	MHTsallis_HR
$I_{R,\alpha} \xrightarrow{\alpha} D_{R,\alpha}$	(8)	MIRenyi_DR
$I_{R,\alpha} \xrightarrow{\alpha} H_{R,\alpha}$	(73)	MIRenyi_HR
$I_{T,\alpha} \xrightarrow{\alpha} D_{T,\alpha}$	(8)	MITsallis_DT
$D_{SB,\alpha} \xrightarrow{\alpha} D_{NB,\alpha}$	(36)	MDSymBregman_DB
$D_{JR,\alpha}^\pi \xrightarrow{\alpha} H_{R,\alpha}$	(42)	MDJR_HR
$D_{JT,\alpha} \xrightarrow{\alpha} H_{T,\alpha}$	(43)	MDJT_HT
$K_{JS} \xrightarrow{\pi=[\frac{1}{2};\frac{1}{2}]} D_{JS}^\pi$	(60)	MKJS_DJS
$K_{JT,\alpha} \xrightarrow{\alpha} H_{T,\alpha}$	(61)	MKJT_HT
$K_{EJS,u} \xrightarrow{\pi=[\frac{1}{2};\frac{1}{2}]} D_{JS}^\pi$	(62)	MKExpJS_DJS
$K_{EJR1,u,\alpha} \xrightarrow{\alpha} H_{R,\alpha}$	(63)	MKExpJR1_HR
$K_{EJR2,u,\alpha} \xrightarrow{\pi=[\frac{1}{2};\frac{1}{2}],\alpha} D_{JR,\alpha}^\pi$	(64)	MKExpJR2_DJR
$K_{EJT1,u,\alpha} \xrightarrow{\alpha} H_{T,\alpha}$	(65)	MKExpJT1_HT
$K_{EJT2,u,\alpha} \xrightarrow{\alpha} D_{JT,\alpha}$	(66)	MKExpJT2_DJT

Table 12: Parameter passing in meta estimators. Notation $X \xrightarrow{z} Y$: the X meta method sets the z parameter(s) of the Y estimator. 2nd column: the equation describing this action. 3rd column: cost name.

B Python ITE \leftrightarrow Matlab ITE

Python cost names with Matlab equivalents (when it exists) are summarized in [ite/ite/cost/x_python_to_matlab.py](#):

- Python \rightarrow Matlab cost name transition: see dictionary `dict_X_PythonToMatlab`, where $X \in \{A, C, D, H, I, K, condH, condI\}$.
- Matlab \rightarrow Python conversion, given a cost type X : see `dict_X_MatlabToPython`.

The equations of the estimators can be looked up by this correspondence and the Matlab ITE documentation.¹⁴

C For Mathy People: Axioms of Concordance and Dependence

This section summarizes the axiomatic formulations of concordance (Def. 1, 2, 3) and dependence (Def. 4, 5).

Definition 1 (concordance ordering) *In two dimensions ($d = 2$) a C_1 copula is said to be smaller than the C_2 copula ($C_1 \prec C_2$) [4], if*

$$C_1(\mathbf{u}) \leq C_2(\mathbf{u}), \quad (\forall \mathbf{u} \in [0, 1]^2). \quad (77)$$

This pointwise partial ordering on the set of copulas is called concordance ordering. In the general ($d \geq 2$) case, a C_1 copula is said to be smaller than the C_2 copula ($C_1 \prec C_2$) [2], if

$$C_1(\mathbf{u}) \leq C_2(\mathbf{u}) \text{ and } \bar{C}_1(\mathbf{u}) \leq \bar{C}_2(\mathbf{u}) \quad (\forall \mathbf{u} \in [0, 1]^d). \quad (78)$$

Note:

- ‘ \prec ’ is called concordance ordering; it again defines a partial ordering.
- The rationale behind requiring $C_1 \leq C_2$ and $\bar{C}_1 \leq \bar{C}_2$ is that we want to capture ‘simultaneously large’ and ‘simultaneously small’ tendencies.
- The two definitions [(77), (78)] coincide only in the two-dimensional ($d = 2$) case.

Definition 2 (measure of concordance [5, 3, 4]) *A κ numeric measure of association on pairs of random variables (y^1, y^2 whose joint copula is C) is called a measure of concordance, if it satisfies the following properties:*

¹⁴Available at <https://bitbucket.org/szzoli/ite/downloads>; see Section E (Estimation Formulas).

A1. Domain: it is defined for every (y^1, y^2) pair of continuous random variables.

A2. Range: $\kappa(y^1, y^2) \in [-1, 1]$, $[\kappa(y^1, y^1) = 1, \text{ and } \kappa(y^1, -y^1) = -1]$.

A3. Symmetry: $\kappa(y^1, y^2) = \kappa(y^2, y^1)$.

A4. Independence: if y^1 and y^2 are independent, then $\kappa(y^1, y^2) = \kappa(\Pi) = 0$.

A5. Change of sign: $\kappa(-y^1, y^2) = -\kappa(y^1, y^2)$ [= $\kappa(y^1, -y^2)$].

A6. Coherence: if $C_1 \prec C_2$, then $\kappa(C_1) \leq \kappa(C_2)$.¹⁵

A7. Continuity: if (y_t^1, y_t^2) is a sequence of continuous random variables with copula C_t , and if C_t converges to C pointwise, then $\lim_{t \rightarrow \infty} \kappa(C_t) = \kappa(C)$.

Note: the properties in brackets can be derived from the others.

Definition 3 (multivariate measure of concordance [1, 7]) A multivariate measure of concordance is a κ function that assigns to every continuous random variable \mathbf{y} a real number and satisfies the following requirements:

B1. Normalization:

B1a : $\kappa(y^1, \dots, y^d) = 1$ if each y^i is an increasing function of every other y^j (or in terms of copulas $\kappa(M) = 1$), and

B1b : $\kappa(y^1, \dots, y^d) = 0$ if y^i -s are independent (or in terms of copulas $\kappa(\Pi) = 1$).

B2. Monotonicity: $C_1 \prec C_2 \Rightarrow \kappa(C_1) \leq \kappa(C_2)$.

B3. Continuity: If the cdf of the random variable sequence \mathbf{y}_t (F_t) converges to F , the cdf of \mathbf{y} ($\lim_{t \rightarrow \infty} F_t = F$), then $\lim_{t \rightarrow \infty} \kappa(\mathbf{y}_t) = \kappa(\mathbf{y})$. [In terms of copulas: $\lim_{t \rightarrow \infty} C_t = C$ (uniformly) $\Rightarrow \lim_{t \rightarrow \infty} \kappa(C_t) = \kappa(C)$.]

B4. Permutation invariance: if $\{i_1, \dots, i_d\}$ is permutation of $\{1, \dots, d\}$, then $\kappa(y^{i_1}, \dots, y^{i_d}) = \kappa(y^1, \dots, y^d)$.

B5. Duality: $\kappa(-y^1, \dots, -y^d) = \kappa(y^1, \dots, y^d)$.

B6. Reflection symmetry property: $\sum_{\epsilon_1, \dots, \epsilon_d = \pm 1} \kappa(\epsilon_1 y^1, \dots, \epsilon_d y^d) = 0$, where the sum is over all the 2^d possibilities.

B7. Transition property: there exists a sequence of r_d numbers such that for all \mathbf{y} $r_{d-1} \kappa(y^2, \dots, y^d) = \kappa(y^1, \dots, y^d) + \kappa(-y^1, \dots, y^d)$.

Definition 4 (measure of dependence) [4] defined a numeric measure κ between two random variables y^1 and y^2 whose copula is C as a measure of dependence if it satisfies the following properties:

C1. Domain: κ is defined for every (y^1, y^2) pair.

C2. Symmetry: $\kappa(y^1, y^2) = \kappa(y^2, y^1)$.

C3. Range: $\kappa(y^1, y^2) \in [0, 1]$.

C4. Independence: $\kappa(y^1, y^2) = 0$ if and only if y^1 and y^2 are independent.

C5. Strictly monotone functional dependence: $\kappa(y^1, y^2) = 1$ if and only if each of y^1 and y^2 is a strictly monotone function of the other.

C6. Invariance to strictly monotone functions: if f_1 and f_2 are strictly monotone functions, then $\kappa(y^1, y^2) = \kappa(f_1(y^1), f_2(y^2))$.

C7. Continuity: if (y_t^1, y_t^2) is a sequence of random variables with copula C_n , and if $\lim_{t \rightarrow \infty} C_t = C$ (pointwise), then $\lim_{t \rightarrow \infty} \kappa(C_t) = \kappa(C)$.

Definition 5 (multivariate measure of dependence) [8] defined the notion of measure of dependence in case of d dimension as follows. A κ real-valued function is called a measure of dependence if it satisfies the properties:

¹⁵Hence the name concordance ordering.

D1. Domain: κ is defined for any continuously distributed \mathbf{y} .

D2. Permutation invariance: if $\{i_1, \dots, i_d\}$ is permutation of $\{1, \dots, d\}$, then $\kappa(y^{i_1}, \dots, y^{i_d}) = \kappa(y^1, \dots, y^d)$.

D3. Normalization: $0 \leq \kappa(y^1, \dots, y^d) \leq 1$.

D4. Independence: $\kappa(y^1, \dots, y^d) = 0$ if and only if y^i -s are independent.

D5. Strictly monotone functional dependence: $\kappa(y^1, \dots, y^d) = 1$ if and only if each y^i is an increasing function of each of the others.

D6. Invariance to strictly monotone functions: If f_1, \dots, f_d are all strictly increasing functions, then $\kappa(y^1, \dots, y^d) = \kappa(f_1(y^1), \dots, f_d(y^d))$.

D7. Normal case: Let \mathbf{y} be normally distributed and $\rho_{ij} = \text{cov}(y^i, y^j)$. If r_{ij} -s are either all non-negative, or all non-positive then κ is a strictly increasing function of each of the $|r_{ij}|$ -s.

D8. Continuity: If the random variable sequence \mathbf{y}_t converges in distribution to \mathbf{y} , then $\lim_{t \rightarrow \infty} \kappa(\mathbf{y}_t) = \kappa(\mathbf{y})$.

References

- [1] Ali Dolati and Manuel Úbeda-Flores. On measures of multivariate concordance. *Journal of Probability and Statistical Science*, 4:147–164, 2006.
- [2] Harry Joe. Multivariate concordance. *Journal of Multivariate Analysis*, 35:12–30, 1990.
- [3] Roger B. Nelsen. *Distributions with Given Marginals and Statistical Modelling*, chapter Concordance and copulas: A survey, pages 169–178. Kluwer Academic Publishers, Dordrecht, 2002.
- [4] Roger B. Nelsen. *An Introduction to Copulas*. Springer, 2006.
- [5] Marco Scarsini. On measures of concordance. *Stochastica*, 8:201–218, 1984.
- [6] Zoltán Szabó. Information theoretical estimators toolbox. *Journal of Machine Learning Research*, 15:283–287, 2014.
- [7] M. D. Taylor. Multivariate measures of concordance. *Annals of the Institute of Statistical Mathematics*, 59:789–806, 2007.
- [8] Edward F. Wolff. N-dimensional measures of dependence. *Stochastica*, 4:175–188, 1980.