

Statistical independence measure based on maximum norm of joint and product-marginal characteristic functions

Povilas Daniušis, povilasd@neurotechnology.com*

November 2021

Abstract

In this paper we propose statistical independence measure based on the maximum norm of difference between joint and product-marginal characteristic functions, and its estimation procedure (including open-source repository). We also extend the proposed measure to the reproducing kernel Hilbert spaces (RKHS), which allows to apply it to structured data.

We conduct experiments both with simulated and real data. The conducted experiments reveal, that the proposed measure can detect statistical dependence in non-linearly related data sets, and that it can improve real data classification accuracy, when applied for feature extraction and regularisation, on a set of common classifiers.

1 Introduction

Statistical dependence measures plays important role in various statistical and machine learning methods (e.g. hypothesis testing [1], feature selection and extraction [2, 3], information bottleneck methods [4], causal inference [5], self-supervised learning [6], among others). Earliest statistical dependence estimation ideas (e.g. conditional probability) share nearly-common origin with the beginning of formal statistical reasoning itself. During last two centuries ideas of correlation and (relative) entropy (including various generalizations) were proposed and became very popular in numerous applications and theoretical developments. However, with the increasing popularity of statistical machine learning, new statistical dependence estimation methods, that are robust, applicable to noisy, high-dimensional, structured data, and which can be efficiently integrated with modern machine learning methods are helpful for the development both of the theory and application.

In this article we focus on quantitative estimation of statistical independence, using characteristic functions. We begin with the short review of some important previous dependence estimation approaches (Section 2), devoting special

*This preprint currently is not an official product of Neurotechnology

attention to ones based on characteristic functions (Section 2.1). Afterwards, we formulate the proposed characteristic function-based statistical dependence measure and its empirical estimator (Section 3), including an extension into reproducing kernel Hilbert spaces (RKHS'es), which are the main theoretical contribution of our paper. Section 4 is devoted to experiments with simulated and real data sets, where we apply the proposed dependence measure for feature extraction and deep neural network (DNN) regularisation, and finalizing Section 5 concludes this article.

2 Previous Work

During recent years, various approaches have been used in order to construct statistical dependence estimation methods. For example, information theory (mutual information [7] and generalisations), reproducing kernel Hilbert spaces (Hilbert-Schmidt independence criterion [1]), characteristic functions (distance correlation [8, 9]), and other (e.g. [10] copula-based kernel dependence measures, integral-probability-metric-reliant Sobolev independence criterion [11]). Further we will focus on characteristic-function based methods.

2.1 Characteristic-function-based methods

Characteristic function of d_X -dimensional random vector X defined in some probability space $(\Omega_X, \Sigma_X, \mathbb{P}_X)$ is defined as

$$\phi_X(\alpha) = \mathbb{E}_X e^{i\alpha^T X}, \quad (1)$$

where $i = \sqrt{-1}$, $\alpha \in R^{d_X}$. Having n i.i.d. realisations of X , corresponding empirical characteristic function is defined as

$$\widehat{\phi}_X(\alpha) = \frac{1}{n} \sum_{j=1}^n e^{i\langle \alpha, x_j \rangle}. \quad (2)$$

Having pair of two random vectors (X, Y) defined in another probability space $(\Omega_{X,Y}, \Sigma_{X,Y}, \mathbb{P}_{X,Y})$ joint characteristic function is defined as:

$$\phi_{X,Y}(\alpha, \beta) = \mathbb{E}_{X,Y} e^{i(\alpha^T X + \beta^T Y)}, \quad (3)$$

where $\alpha \in R^{d_X}$ and $\beta \in R^{d_Y}$. Similarly, having n i.i.d. realisations of (X, Y) , joint empirical characteristic function is defined as

$$\widehat{\phi}_{X,Y}(\alpha, \beta) = \frac{1}{n} \sum_{j=1}^n e^{i(\langle \alpha, x_j \rangle + \langle \beta, y_j \rangle)}. \quad (4)$$

In terms of characteristic functions, statistical independence of X and Y is equivalent to $\forall \alpha \in R^{d_X}, \forall \beta \in R^{d_Y}$,

$$\Delta_{X,Y}(\alpha, \beta) := \phi_{X,Y}(\alpha, \beta) - \phi_X(\alpha)\phi_Y(\beta) = 0, \quad (5)$$

where d_x and d_y are dimensions of X and Y , respectively.

This formulation of statistical independence was used as the basis (first in [8] for one-dimensional case, and afterwards extended and developed by [9] for bivariate multidimensional random vectors) for construction of statistical independence tests and measures. *Distance covariance* and *distance correlation*, proposed by [9] relies on weighted L^2 -norm analysis of (5). They select weighting function in such a way, that dependence measure can be expressed in terms of correlation of data-dependent distances. Study [12] generalises [9] to multivariable case and proposes *distance multivariate* and derivative dependence measure, called *total distance multivariate*.

Our motivation stems from the fact that evaluation of [9] measures in high dimensional cases may be prone to curse of dimensionality (as mentioned in [13]). Although staying in L^p -space framework, instead of $p = 2$ (L^2 space) we take a limit when $p \rightarrow \infty$, and also avoid direct calculation of integral by working in L^∞ , which has corresponding supremum norm. Also, from practical point of view maximization is convenient, because it is already implemented in various deep learning frameworks.

3 Proposed Independence Measure

This motivates the construction of a novel dependence measure, which we further refer to as Kac independence measure (KacIM). Let X and Y be two standartized random random vectors. The proposed independence measure is defined as

$$\kappa(X, Y) = \max_{\alpha \in \mathbb{R}^{d_X}, \beta \in \mathbb{R}^{d_Y}} |\phi_{X,Y}(\alpha, \beta) - \phi_X(\alpha)\phi_Y(\beta)|. \quad (6)$$

3.1 Properties

Theorem 1. *Statistical independence measure (6) has the following properties:*

1. $\kappa(X, Y) = \kappa(Y, X)$,
2. $0 \leq \kappa(X, Y) \leq 1$,
3. $\kappa(X, Y) = 0$ iff $X \perp Y$.
4. $\kappa(X, Y)$ is scale invariant.

Proof. Property 1. is obvious from definition (6) (commutativity of addition and multiplication), and property 2. directly follows from Cauchy inequality and that absolute value of characteristic function is bounded by 1:

$$\begin{aligned} |\phi_{X,Y}(\alpha, \beta) - \phi_X(\alpha)\phi_Y(\beta)|^2 &= \mathbb{E}_{X,Y} |(e^{i\alpha^T X} - \phi_X(\alpha))(e^{i\beta^T Y} - \phi_Y(\beta))|^2 \leq \\ \mathbb{E}_{X,Y} |e^{i\alpha^T X} - \phi_X(\alpha)|^2 |e^{i\beta^T Y} - \phi_Y(\beta)|^2 &= (1 - |\phi_X(\alpha)|^2)(1 - |\phi_Y(\beta)|^2). \end{aligned}$$

Proof of property 3. directly follows from properties of characteristic functions (see e.g. [14], Corollary 14.1)¹. Scale invariance (Property 4.) is trivial result of the standartisation requirement for X and Y . \square

3.2 Estimation

Having i.i.d. standartized data (x_j, y_j) , $j = 1, 2, \dots, n$, an empirical scale-invariant estimator of (6) is defined via corresponding empirical characteristic functions (4) and (2):

$$\hat{\kappa}(X, Y) = \max_{\alpha, \beta} |\widehat{\phi}_{X,Y}(\alpha, \beta) - \widehat{\phi}_X(\alpha) \widehat{\phi}_Y(\beta)|. \quad (7)$$

Empirical estimator (7) also is symmetric and and bounded (Theorem 1). Empirically we observed that normalisation of parameters α and β on to unit sphere increases estimation stability (why?). The estimator (7) can be calculated iteratively by Algorithm 1 (Pytorch [16] implementation can be accessed from https://github.com/povidanius/kac_independence_measure).

Algorithm 1 KacIM estimation iteration

Require: data batch (x, y) , gradient-based optimiser $GradOpt(loss)$
 Normalize (x, y) to zero mean and unit variance (scale invariance).
 Calculate KacIM estimator $\hat{\kappa}(x, y)$, without maximization step (i.e. using current α, β).
 Perform one maximization iteration of computed $\hat{\kappa}(x, y)$ via $\alpha, \beta \rightarrow GradOpt(\hat{\kappa}(x, y))$.

In our implementation we use decoupled weight decay regularization optimizer [17].

3.3 Kernel version

Having two RKHS'es, defined by feature mapping $(x, y) \rightarrow (l(x, \cdot), l(y, \cdot))$, where $k : \mathbb{R}^{d_x} \times \mathbb{R}^{d_x} \rightarrow \mathbb{R}$ and $l : \mathbb{R}^{d_y} \times \mathbb{R}^{d_y} \rightarrow \mathbb{R}$ (see [18]). Then, estimation of kernel- KIM ($\hat{\kappa}_{k,l}(X, Y)$) can be reformulated as maximization of :

$$\left| \frac{1}{n} \sum_{j=1}^n e^{i \langle \alpha, k(x_j, \cdot) \rangle + \langle \beta, l(y_j, \cdot) \rangle} - \frac{1}{n^2} \sum_{j=1}^n e^{i \langle \alpha, k(x_j, \cdot) \rangle} \sum_{k=1}^n e^{i \langle \beta, l(y_k, \cdot) \rangle} \right|, \quad (8)$$

and representer theorem[?] implies

$$\hat{\kappa}_{k,l}(X, Y) = \max_{\|\alpha\|=\|\beta\|=1} \left| \frac{1}{n} 1^T e^{i(\alpha^T K + \beta^T L)} - \frac{1}{n^2} (1^T e^{i(\alpha^T K)}) (1^T e^{i(\beta^T L)}) \right|, \quad (9)$$

¹This property also is known as Kac's theorem [15]. Although it is quite simple mathematical fact, this provides the basis of the proposed measure name.

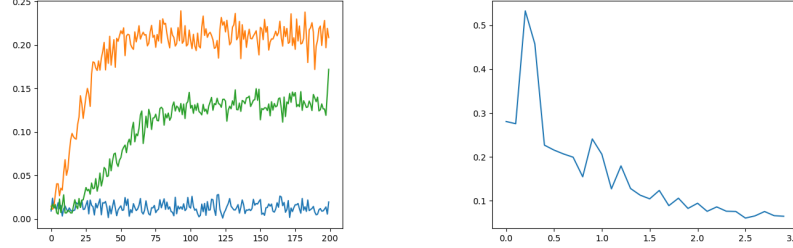


Figure 1: Left figure: Dependence detection in independent data (blue), additive (orange) and multiplicative (green) noise scenarios. Right figure: noise level (x axis) vs final iteration KacIM value (y axis). KacIM values for larger noise levels saturates as in tail of graph

where K and L are Gram matrices, corresponding to x_i and y_i . Note that the number of parameters of $\hat{\kappa}_{k,l}(X, Y)$ is dimension-independent and is equal to $2n_b$, where n_b is batch size. Also, kernel- KIM can be applied for structured data, via corresponding positive defined kernels.

4 Experiments

Further we will conduct empirical investigation of KacIM. We will begin with simple illustrative simulations, and afterwards investigate empirical performance of KacIM in classifier feature extraction and regularization tasks using various publicly available data sets.

4.1 Generated data

Non-linear statistical dependence detection. We begin with simple example, which demonstrates the efficiency of KacIM for simulated multivariate data with additive and multiplicative noise.

Figure 1 reflects KacIM values during iterative adaptation (200 iterations). In the case of independent data, both x_i and y_i ($d_x = 512$, $d_y = 4$) are sampled from gaussian distribution, independently. In the case of dependent data, an additive noise and multiplicative noise, the dependent variable is generated according to $y_i = \sin(Px_i) + \cos(Px_i) + \lambda\epsilon_i$ ($\lambda = 1.00$) and $y_i = (\sin(Px_i) + \cos(Px_i))\epsilon_i$, respectively, where P is $d_x \times d_y$ random projection matrix, $\epsilon_i \sim N(0, 1)$ and $\epsilon_i \perp x_i$.

When data is independent, both in additive and multiplicative cases, due to independence, estimator (7) is resistant to maximisation, and oscillates near zero. On the other hand, when the data is not independent, the condition (5) is violated and maximization of estimator (7) is possible.

Noise variance effect In this simulation we use the same additive noise setting as in previous paragraph, but evaluate all noise levels $\lambda \in [0.1, 3.0]$, with step 0.1. Figure 1 empirically shows that value of KacIM negatively correlates with noise level, and therefore the proposed measure is able not only to detect whether independence is present, but also to quantitatively evaluate it, which enables to use it to derive cost functions for various learning-based algorithms.

4.2 Feature selection

4.3 Feature extraction

In feature extraction experiments we will use a set of classification data sets from OpenML [19]. We conduct linear feature extraction by seeking

$$W^* = \arg \max_W \kappa(Wx, y) - \alpha \text{Tr}\{(W^T W - I)^T (W^T W - I)\}, \quad (10)$$

where the regularisation term, controlled by multiplier $\alpha \geq 0$, enforces projection matrix W^* to be orthogonal.

Afterwards, feature extraction is conducted by $f = W^*x$ and these features are used as the inputs to logistic regression classifier.

We randomly split all the datasets in training and testing sets of equal size, comparing unmodified inputs x , and features of all possible dimensions up to d_x with 10% step. In our experiments we set α to 1.0 to quickly ensure orthogonal projection matrices, and further proceed to dependence maximization stage. The classification accuracies, reported in Table 1 demonstrate that this KacIM-based feature extraction procedure indeed allows to increase classification accuracy when applied to real data sets from various domains.

The purpose of these experiments is to provide the preliminary evaluation of the applicability of KacIM for feature extraction, hence we use rather basic comparative baselines.

4.4 Regularisation

In regularisation experiments we investigate chest x-ray classification task. It is represented as binary classification data set, consisting of xray scans, which should be classified as pneumonia or normal.

Let $f(x|\theta)$ be DNN classifier. We will investigate additive regularizer, which maximises dependency of bottleneck feature $\phi(x)$ and target variable y :

$$\text{Cost}(\theta) = CE(f(x|\theta), y) + \beta \kappa(\phi(x|\theta), y), \quad (11)$$

where CE is cross-entropy loss, and $\beta \geq 0$ is regularisation parameters.

5 Conclusion

In this article we propose statistical dependence measure, KacIM, which relies on the L^∞ norm of the absolute value of difference between joint characteristic

Classification accuracies				
Dataset	size/dim/ n_c .	Raw	KacIMFE	NCA
micro-mass	(360,1300,10)	0.88	0.93	0.89
ionosphere	(351,34,2)	0.88	0.92	0.94
spambase	(4601,57,2)	0.69	0.83	0.87
lsvt	(126,310,2)	0.67	0.68	0.78
one-hundred-plants-texture	(1599,64,100)	0.62	0.77	0.80
tokyo1	(959,44,2)	0.72	0.89	0.91
clean1	(476,168,2)	0.77	0.99	0.98
one-hundred-plants-shape	(1600,64,100)	0.10	0.48	0.47
pc4	(1458,37,2)	0.88	0.87	0.87
madelon	(2600,500,2)	0.60	0.55	0.57
scene	(2407,299,2)	0.89	0.98	0.93
gina_agnostic	(3468,970,2)	0.85	0.80	0.81
isolet	(7797,617,26)	0.93	0.95	0.95

Table 1: Classification accuracies, trained on KacIM-based features.

function and the product of marginal ones. The proposed measure, in theory can detect arbitrary statistical dependence between a pairs of random variables of different dimension, extended to various known statistical generalisations (e.g. reproducing kernel Hilbert spaces, multivariablity), machine learning scenarios (e.g. feature extraction, regularisation, among others), and is empirically tractable on these problems. On the other side, it raises a corresponding set of unanswered questions, both theoretical and empirical. For example, the interpretability when this measure approaches its maximal value remains insufficiently clear, however empirical experiments with simulated data reveals, that increasing independence between two random variables is reflected in a decreasing trend on the estimated values of the proposed dependence measure. Beside direct applications, the proposed measure is differentiable and thereby can be integrated with modern deep-learning methods, applied to high-dimensional and structured data. From empirical point of view, we see exploration of KacIM in causality, information bottleneck, self-supervised learning, and other modern problems, where dependence measures define a criterion of optimisation as important future work.

6 Acknowledgements

References

- [1] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *ALT*, 2005.

- [2] P. Daniušis and P. Vaitkus. Supervised feature extraction using hilbert-schmidt norms. In Emilio Corchado and Hujun Yin, editors, *Intelligent Data Engineering and Automated Learning - IDEAL 2009*, pages 25–33, Berlin, Heidelberg, 2009. Springer Berlin Heidelberg.
- [3] P. Daniušis, Pr. Vaitkus, and L. Petkevičius. Hilbert-schmidt component analysis. *Lithuanian mathematical journal*, 57(A):7–11, Dec. 2016.
- [4] Kurt Wan-Duo Ma, J. P. Lewis, and W. Kleijn. The hsic bottleneck: Deep learning without back-propagation. *ArXiv*, abs/1908.01580, 2020.
- [5] Patrik Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2009.
- [6] Yazhe Li, Roman Pogodin, Danica J. Sutherland, and Arthur Gretton. Self-supervised learning with kernel dependence maximization. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [7] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory 2nd Edition (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, July 2006.
- [8] Andrey Feuerverger. A consistent test for bivariate dependence. *International Statistical Review / Revue Internationale de Statistique*, 61(3):419–433, 1993.
- [9] Gábor J. Székely, Maria L. Rizzo, and Nail K. Bakirov. Measuring and testing dependence by correlation of distances. *The Annals of Statistics*, 35(6):2769 – 2794, 2007.
- [10] Barnabás Póczos, Zoubin Ghahramani, and Jeff G. Schneider. Copula-based kernel dependency measures. *ArXiv*, abs/1206.4682, 2012.
- [11] Youssef Mroueh, Tom Sercu, Mattia Rigotti, Inkit Padhi, and Cicero Nogueira dos Santos. Sobolev independence criterion. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 9505–9515. Curran Associates, Inc., 2019.
- [12] Björn Böttcher, Martin Keller-Ressel, and René Schilling. Distance multivariate: New dependence measures for random vectors, 10 2018.
- [13] Dominic Edelmann, Konstantinos Fokianos, and Maria Pitsillou. An Updated Literature Review of Distance Correlation and Its Applications to Time Series. *International Statistical Review*, 87(2):237–262, August 2019.
- [14] Jacod Jean. *Probability essentials / Jean Jacod, Philip Protter*. Universitext. Springer, Berlin Heidelberg New York, 2nd edition edition, 2003.
- [15] Johan Kustermans J. Martin Lindsay Michael Schuermann Uwe Franz David Applebaum, B.V. Rajarama Bhat. Quantum independent increment processes i: From classical probability to quantum stochastic calculus. 2005.
- [16] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [17] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- [18] Bernhard Schölkopf, Alexander J. Smola, and Francis Bach. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. The MIT Press, 2018.
- [19] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.