

Exploring and measuring non-linear correlations: Copulas, Lightspeed Transportation and Clustering

Gautier Marti
Hellebore Capital Ltd
Ecole Polytechnique

Sébastien Andler
ENS de Lyon
Hellebore Capital Ltd

Frank Nielsen
Ecole Polytechnique
LIX - UMR 7161

Philippe Donnat
Hellebore Capital Ltd
Michelin House, London

Abstract

We propose a methodology to explore and measure the pairwise correlations that exist between variables in a dataset. The methodology leverages copulas for encoding dependence between two variables, state-of-the-art optimal transport for providing a relevant geometry to the copulas, and clustering for summarizing the main dependence patterns found between the variables. Some of the clusters centers can be used to parameterize a novel dependence coefficient which can target or forget specific dependence patterns. Finally, we illustrate and benchmark the methodology on several datasets. Code and numerical experiments are available online for reproducible research.

Introduction

Pearson's correlation coefficient which estimates linear dependence between two variables is still the mainstream tool for measuring variable correlations in science and engineering. However, its shortcomings are well-documented in the statistics literature: not robust to outliers; not invariant to monotone transformations of the variables; can take value 0 whereas variables are strongly dependent; only relevant when variables are jointly normally distributed. A large but under-exploited literature in statistics and machine learning has expanded recently to alleviate these issues (Reshef et al. 2011; Székely, Rizzo, and others 2009; Sejdinovic et al. 2013). An underlying idea to many of the dependence coefficients is to compute a distance $D(P(X, Y), P(X)P(Y))$ between the joint distribution $P(X, Y)$ of variables X, Y and $P(X)P(Y)$ the product of marginal distributions encoding the independence. For example, choosing $D = \text{KL}$ (Kullback-Leibler divergence), we end up with the Mutual Information (MI) measure, well-known in information theory. Thus, one can detect all the dependences between X and Y since the distance will be greater than 0 as soon as $P(X, Y)$ is different from $P(X)P(Y)$. Then, the dependence literature focus has shifted toward the new concept of "equitability" (Kinney and Atwal 2014): How can one quantify the strength of a statistical association between two variables without bias for relationships of a specific form? Many researchers now aim at designing and proving that their proposed measures are indeed equitable (Reshef et al. 2013; Ding and Li 2013; Chang et al. 2016). This is *not* what we look for in this article. But, on the contrary, we want to

target specific dependence patterns and ignore others. We want to target dependence which are relevant to such or such problem, and forget about the dependence which are not in the scope of the problems at hand, or even worse which may be spurious associations (pure chance or artifacts in the data). The latter will be detected with an equitable dependence measure since they are deviation from independence, and will be given as much weight as the interesting ones. Rather than using the biases for specific dependence of several coefficients, we propose a dependence coefficient that can be parameterized by a set of *target-dependences*, and a set of *forget-dependences*. Sets of target and forget dependences can be built using expert hypotheses, or by leveraging the centers of clusters resulting from an exploratory clustering of the pairwise dependences. To achieve this goal, we will leverage three tools: copulas, optimal transportation, and clustering. Whereas clustering, the task of grouping a set of objects in such a way that objects in the same group (also called cluster) are more similar to each other than those in different groups, is common knowledge in the machine learning community, copulas and optimal transportation are not yet mainstream tools. Copulas have recently gained attention in machine learning (Elidan 2013), and several copula-based dependence measures have been proposed for improving feature selection methods (Ghahramani, Póczos, and Schneider 2012; Lopez-Paz, Hennig, and Schölkopf 2013; Chang et al. 2016). Optimal transport may be more familiar to computer scientists working in computer vision since it is the underlying theory of the Earth Mover's Distance (Rubner, Tomasi, and Guibas 2000). Until very recently, optimal transportation distances between distributions were not deemed relevant for machine learning applications since the best computational cost known was super-cubic to the number of bins used for discretizing the distribution supports which grows itself exponentially with the dimension. A mere distance evaluation could take several seconds! In this article, we leverage recent computational breakthroughs detailed in (Cuturi 2013) which make their use practical in machine learning.

Background on Copulas and Optimal Transport

Copulas

Copulas are functions that couple multivariate distribution functions to their one-dimensional marginal distribution functions (Nelsen 2013). In this article, we will only consider bivariate copulas, but most of the results and the methodology presented hold in the multivariate setting, at the cost of a much higher computational burden which is for now a bit unrealistic.

Theorem 1 (Sklar’s Theorem (Sklar 1959)) *For any random vector $X = (X_i, X_j)$ having continuous marginal cumulative distribution functions F_i, F_j respectively, its joint cumulative distribution function F is uniquely expressed as $F(X_i, X_j) = C(F_i(X_i), F_j(X_j))$, where C , the bivariate distribution of uniform marginals $U_i, U_j := F_i(X_i), F_j(X_j)$, is known as the copula of X .*

Copulas are central for studying the dependence between random variables: their uniform marginals jointly encode all the dependence. They allow to study scale-free measures of dependence and are *invariant to monotonous transformations of the variables*. Some copulas play a major role in the measure of dependence, namely \mathcal{W} and \mathcal{M} the Fréchet-Hoeffding copula bounds, and the independence copula $\Pi(u_i, u_j) = u_i u_j$ (depicted in Figure 1).

Definition 1 (Fréchet-Hoeffding copula bounds) *For any copula $C : [0, 1]^2 \rightarrow [0, 1]$ and any $(u_i, u_j) \in [0, 1]^2$ the following bounds hold:*

$$\mathcal{W}(u_i, u_j) \leq C(u_i, u_j) \leq \mathcal{M}(u_i, u_j), \quad (1)$$

where $\mathcal{W}(u_i, u_j) = \max\{u_i + u_j - 1, 0\}$ is the copula for countermonotonic random variables and $\mathcal{M}(u_i, u_j) = \min\{u_i, u_j\}$ is the copula for comonotonic random variables.

Many correlation coefficients can actually be expressed as a distance between the data copula and one of these reference copulas. For example, the Spearman (rank) correlation ρ_S which is usually understood as $\rho_S(X_i, X_j) = \rho(F_i(X_i), F_j(X_j))$, i.e. the linear dependence of the probability integral transformed variables (rank-transformed data), can also be viewed as an average distance between the copula C of (X_i, X_j) and the independence copula Π : $\rho_S(X_i, X_j) = 12 \int \int_{[0,1]^2} (C(u_i, u_j) - u_i u_j) du_i du_j$ (Nelsen 2013). Moreover, since $|u_i - u_j|/\sqrt{2}$ is the distance between point (u_i, u_j) to the diagonal (the measure of the positive dependence copula), one can rewrite $\rho_S(X_i, X_j) = 12 \int \int_{[0,1]^2} (C(u_i, u_j) - u_i u_j) du_i du_j = 12 \int \int_{[0,1]^2} u_i u_j dC(u_i, u_j) - 3 = 1 - 6 \int \int_{[0,1]^2} (u_i - u_j)^2 dC(u_i, u_j)$ (Liebscher and others 2014). Thus, Spearman correlation can also be viewed as measuring a deviation from the monotonically increasing dependence to the data copula using a quadratic distance. *We will leverage this idea to propose our dependence-parameterized dependence coefficient.*

Notice that when working with empirical data, we do not know a priori the margins F_i for applying the probability

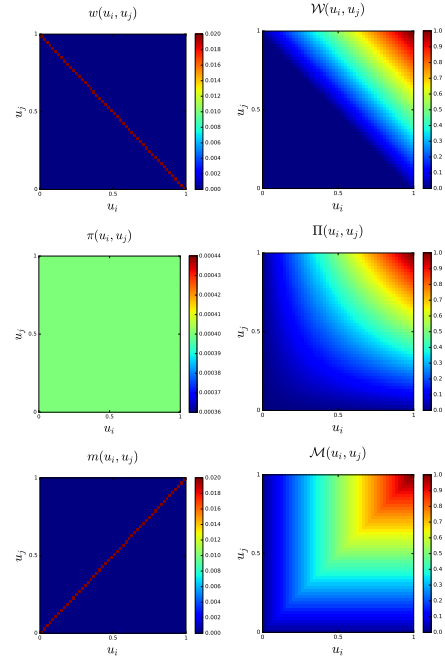


Figure 1: Copulas measure (left column) and cumulative distribution function (right column) heatmaps for negative dependence (first row), independence (second row), i.e. the uniform distribution over $[0, 1]^2$, and positive dependence (third row)

integral transform $U_i := F_i(X_i)$. Deheuvels in (Deheuvels 1979) has introduced a practical estimator for the uniform margins and the underlying copula, the empirical copula transform.

Definition 2 (Empirical Copula Transform) *Let (X_i^t, X_j^t) , $t = 1, \dots, T$, be T observations from a random vector (X_i, X_j) with continuous margins. Since one cannot directly obtain the corresponding copula observations $(U_i^t, U_j^t) := (F_i(X_i^t), F_j(X_j^t))$, where $t = 1, \dots, T$, without knowing a priori F_i , one can instead estimate the empirical margins $F_i^T(x) = \frac{1}{T} \sum_{t=1}^T \mathbf{1}(X_i^t \leq x)$, to obtain the T empirical observations $(\tilde{U}_i^t, \tilde{U}_j^t) := (F_i^T(X_i^t), F_j^T(X_j^t))$. Equivalently, since $\tilde{U}_i^t = R_i^t/T$, R_i^t being the rank of observation X_i^t , the empirical copula transform can be considered as the normalized rank transform.*

Notice that the empirical copula transform is fast to compute, sorting arrays of length T can be done in $O(T \log T)$, consistent and converges fast to the underlying copula (Deheuvels 1981), (Ghahramani, Póczos, and Schneider 2012).

As motivated in the introduction, we want to compare and summarize the pairwise empirical dependence structure (empirical bivariate copulas) of many variables. This brings the following questions: How can we compare two such copulas? What is a relevant representative of a set of empirical copulas? Which geometries are relevant for clustering these empirical distributions, and which are not?

Optimal Transport

In (Marti et al. 2016), authors illustrate in a parametric setting using Gaussian copulas that common divergences (such as Kullback-Leibler, Jeffreys, Hellinger, Bhattacharyya) are not relevant for clustering these distributions, especially when dependence is high. These information divergences are only defined for absolutely continuous measures whereas some copulas have no density (e.g. the one for positive dependence). In practice, when working with frequency histograms, it gets worse: One has to pre-process the empirical measures with a kernel density estimator before computing these divergences. On the contrary, optimal transport distances are well-defined for both discrete (e.g. empirical) and continuous measures.

The idea of optimal transport is intuitive. It was first formulated by Gaspard Monge in 1781 (Monge 1781) as a problem to efficiently level the ground: Given that work is measured by the distance multiplied by the amount of dirt displaced, what is the minimum amount of work required to level the ground? Optimal transport plans and distances give the answer to this problem.

In practice, empirical distributions can be represented by histograms. We follow notations from (Cuturi 2013). Let r, c be two histograms in the probability simplex $\Sigma_m = \{x \in \mathbb{R}_+^m : x^\top \mathbf{1}_m = 1\}$. Let $U(r, c) = \{P \in \mathbb{R}_+^{m \times m} \mid P \mathbf{1}_m = r, P^\top \mathbf{1}_m = c\}$ be the transportation polytope of r and c , that is the set containing all possible transport plans between r and c .

Definition 3 (Optimal Transport) *Given a $m \times m$ cost matrix M , the cost of mapping r to c using a transportation matrix P can be quantified as $\langle P, M \rangle_F$, where $\langle \cdot, \cdot \rangle_F$ is the Frobenius dot-product. The optimal transport between r and c given transportation cost M is thus:*

$$d_M(r, c) := \min_{P \in U(r, c)} \langle P, M \rangle_F. \quad (2)$$

Whenever M belongs to the cone of distance matrices, the optimum of the transportation problem $d_M(r, c)$ is itself a distance.

Lightspeed transportation. Optimal transport distances suffer from a computational burden scaling in $O(m^3 \log m)$ which has prevented their widespread use in machine learning: A mere distance computation between two high-dimensional histograms can take several seconds. In (Cuturi 2013), Cuturi provides a solution to this problem: He restrains the polytope $U(r, c)$ of all possible transport plans between r and c to a Kullback-Leibler ball $U_\alpha(r, c) \subset U(r, c)$, where $U_\alpha(r, c) = \{P \in U(r, c) \mid \text{KL}(P \| rc^\top) \leq \alpha\}$. He then shows that it amounts to perform an entropic regularization (recently generalized to many more regularizers in (Muzellec et al. 2016; Dessein, Papadakis, and Rouas 2016)) of the optimal transportation problem whose solution is smoother and less deterministic. The regularized optimal transportation problem is now strictly convex, and can be solved efficiently using the Sinkhorn-Knopp iterative algorithm which exhibits linear convergence. Its solution is the Sinkhorn distance (Cuturi 2013):

$$d_{M, \alpha}(r, c) := \min_{P \in U_\alpha(r, c)} \langle P, M \rangle_F, \quad (3)$$

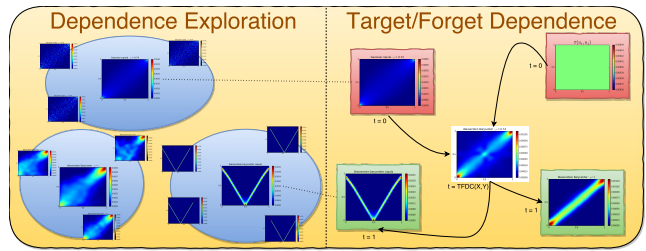


Figure 2: Exploration (left panel) and measure (right panel) of non-linear correlations. Exploration consists in finding clusters of similar copulas, visualizing their centroids, and eventually using them to assess the dependence of given variables represented by their copula

and its dual $d_M^\lambda(r, c): \forall \alpha > 0, \exists \lambda > 0$,

$$d_{M, \alpha}(r, c) = d_M^\lambda(r, c) := \langle P^\lambda, M \rangle_F, \quad (4)$$

where $P^\lambda = \operatorname{argmin}_{P \in U(r, c)} \langle P, M \rangle_F - \frac{1}{\lambda} h(P)$, and h is the entropy function.

In the following, we will leverage the dual-Sinkhorn distances for comparing, clustering and computing the clusters centers (Cuturi and Doucet 2014) of a set of copulas at full speed.

A methodology to explore and measure non-linear correlations

We propose an approach to explore and measure non-linear correlations between N variables X_1, \dots, X_N in a dataset. These N variables can be, for instance, time series or features. The methodology presented (which is summarized in Figure 2) is twofold, and consists of: (i) an exploratory part of the pairwise dependence between variables, (ii) the parameterization and use of a novel dependence coefficient.

Using transportation of copulas as a measure of correlations

In this section, we leverage and extend the idea presented in our short introduction to copulas: correlation coefficients can be viewed as a distance between the data-copula and the Fréchet-Hoeffding bounds or the independence copula. The distance involved is usually an ℓ_p Minkowski metric distance. In the following, we will:

- replace the ℓ_p distance by an optimal transport distance between measures,
- parameterize a dependence coefficient with other copulas than the Fréchet-Hoeffding bounds or the independence one.

Using the optimal transport distance between copulas, we now propose a dependence coefficient which is parameterized by two sets of copulas: *target* copulas and *forget* copulas.

Definition 4 (Target/Forget Dependence Coefficient)

Let $\{C_l^-\}_l$ be the set of forget-dependence copulas. Let

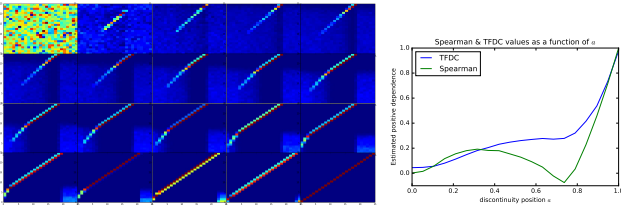


Figure 3: Empirical copulas for (X, Y) where $X = Z\mathbf{1}_{Z < a} + \epsilon_X \mathbf{1}_{Z > a}$, $Y = Z\mathbf{1}_{Z < a+0.25} + \epsilon_Y \mathbf{1}_{Z > a+0.25}$, $a = 0, 0.05, \dots, 0.95, 1$, and where Z is uniform on $[0, 1]$ and ϵ_X, ϵ_Y are independent noises (left figure). Top left is an empirical copula for independence ($a = 0$), bottom right is the copula for perfect positive dependence ($a = 1$). Parameter a is increasing from top to bottom, and from left to right; TFDC and Spearman coefficients estimated between X and Y as a function of a (right figure). For $a = 0.75$, Spearman coefficient yields a negative value, yet $X = Y$ over $[0, a]$

$\{C_k^+\}_k$ be the set of target-dependence copulas. Let C be the copula of (X_i, X_j) . Let d_M be an optimal transport distance parameterized by a ground metric M . We define the Target/Forget Dependence Coefficient as:

$$\text{TFDC}(X_i, X_j; \{C_k^+\}_k, \{C_l^-\}_l) := \frac{\min_l d_M(C_l^-, C)}{\min_l d_M(C_l^-, C) + \min_k d_M(C, C_k^+)} \in [0, 1]. \quad (5)$$

Using this definition, we obtain: $\text{TFDC}(X_i, X_j; \{C_k^+\}_k, \{C_l^-\}_l) = 0 \Leftrightarrow C \in \{C_l^-\}_l$, $\text{TFDC}(X_i, X_j; \{C_k^+\}_k, \{C_l^-\}_l) = 1 \Leftrightarrow C \in \{C_k^+\}_k$.

Example. A standard correlation coefficient can be obtained by setting the forget-dependence set to the independence copula, and the target-dependence set to the Fréchet-Hoeffding bounds. How does it compare to the Spearman correlation? In Figure 3, we display how the two coefficients behave on a simple numerical experiment: $X = Z\mathbf{1}_{Z < a} + \epsilon_X \mathbf{1}_{Z > a}$, $Y = Z\mathbf{1}_{Z < a+0.25} + \epsilon_Y \mathbf{1}_{Z > a+0.25}$, where Z is uniform on $[0, 1]$ and ϵ_X, ϵ_Y are independent noises. That is $X = Y$ over $[0, a]$. Notice that for $a = 0.75$, Spearman coefficient takes a negative value. We may thus prefer the monotonically increasing behaviour of the TFDC to the Spearman one.

How to choose, design and build targets?

We now propose two alternatives for choosing, designing and building the *target* and *forget* copulas: an exploratory data-driven approach and an hypotheses testing approach.

Data-driven: Clustering of copulas Assume we have N variables X_1, \dots, X_N , and T observations for each of them. First, we compute $\binom{N}{2} = O(N^2)$ empirical copulas which represent the dependence structure between all the couples (X_i, X_j) . Then, we summarize all these distributions using a center-based clustering algorithm, and extract the clusters centers using a fast computation of Wasserstein barycenters (Cuturi and Doucet 2014). A given center represents the mean dependence between the couples

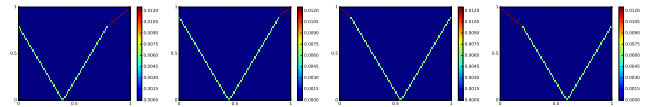


Figure 4: 4 copulas describing the dependence between $X \sim \mathcal{U}([0, 1])$ and $Y \sim (X \pm \epsilon_i)^2$, where ϵ_i is a constant noise specific for each distribution. X and Y are counter-monotonic (more or less) half of the time, and co-monotonic (more or less) half of the time

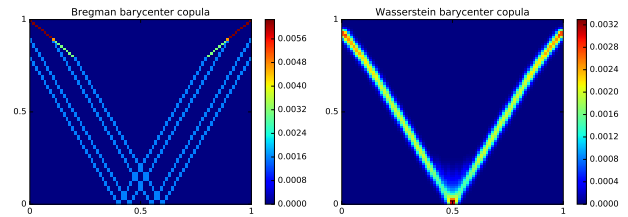


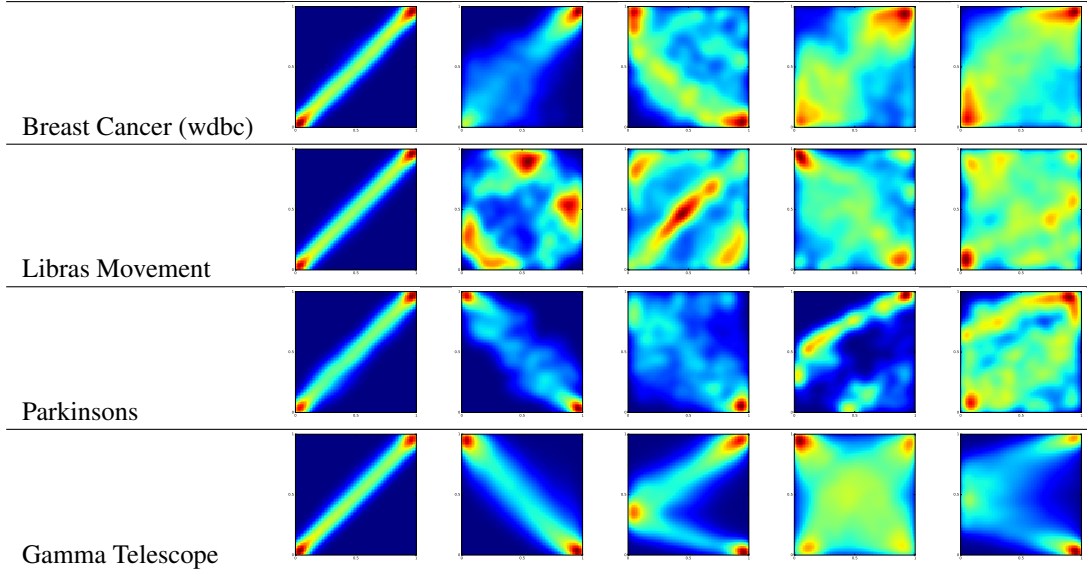
Figure 5: Barycenter of the 4 copulas from Figure 4 for: (left) Bregman geometry (Banerjee et al. 2005) (which includes, for example, squared Euclidean and Kullback-Leibler distances); (right) Wasserstein geometry. Notice that the Wasserstein barycenter better describes the underlying dependence between X and Y : the copula encodes a functional association. This is not the case for the Bregman barycenter

(X_i, X_j) inside the corresponding cluster. Figure 4 and 5 illustrate why a Wasserstein W_2 barycenter, i.e. the minimizer μ^* of $\frac{1}{N} \sum_{i=1}^N W_2^2(\mu, \nu_i)$ (Agueh and Carlier 2011) where $\{\nu_1, \dots, \nu_N\}$ is a set of N measures (here, bivariate empirical copulas), is more relevant to our needs: we benefit from robustness against small deformations of the dependence patterns.

Example. In Table 1, we display some interesting dependence patterns which can be found in UCI datasets <http://archive.ics.uci.edu/ml/>. In this case, variables X_1, \dots, X_N are the N features. Some associations are easy to explain (e.g. top left copula representing the relation between radius and area of roughly round cells in the Breast Cancer Wisconsin (Diagnostic) Data Set) whereas some others less (e.g. top row third copula from the left which represents the relation between the perimeter and the fractal dimension of the cells).

An equitable copula-based dependence measure such as the one described in (Ghahramani, Póczos, and Schneider 2012) may detect them well, but will also detect the spurious ones which are due to artifacts in the data (or pure chance). With this approach, one can spot them and add them to the set of forget-dependence copulas. For these reasons, we think that this approach could improve the feature selection correlation-based approaches (Hall 2000; Yu and Liu 2003) which rely on the hypothesis that *good feature subsets contain features highly correlated with the class, yet uncorrelated with each other* (Hall 2000).

Table 1: Dependence patterns (= clustering centroids) found between variables in UCI datasets



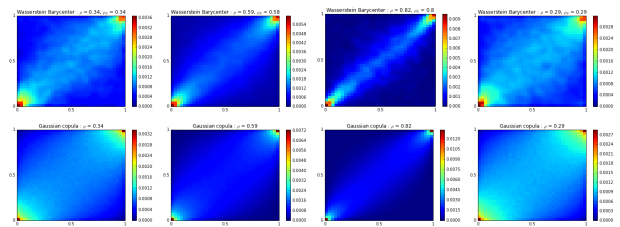
Targets as hypotheses from an expert One can specify dependence hypotheses, generate the corresponding copulas, then measure and rank correlations with respect to them. For example, one can answer to questions such as: Which are the pairs of assets that are usually positively correlated for small variations but uncorrelated otherwise? In (Durante, Saminger-Platz, and Sarkoci 2009), authors present a method for constructing bivariate copulas by changing the values that a given copula assumes on some subrectangles of the unit square. They discuss some applications of their methodology including the construction of copulas with different tail dependencies. Building *target* and *forget* copulas is another one. In the Experiments section, we illustrate its use to answer the previous question and other dependence queries.

Experiments

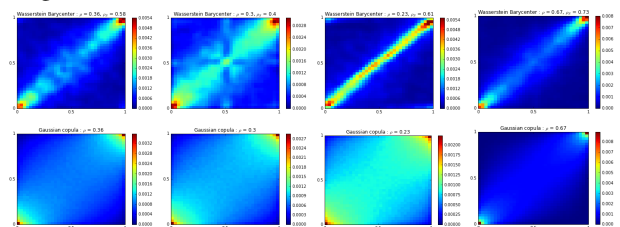
Exploration of financial correlations

We illustrate the first part of the methodology with three different datasets of financial time series. These time series consist in the daily returns of stocks (40 stocks from the CAC 40 index comprising the French highest market capitalizations), credit default swaps (75 CDS from the iTraxx Crossover index comprising the most liquid sub-investment grade European entities) and foreign exchange rates (80 FX rates of major world currencies) between January 2006 and August 2016. We display some of the clustering centroids obtained for each asset class on the top row, and below we display their corresponding Gaussian copulas parameterized by the estimated linear correlations. Notice the strong difference between the empirical copulas and the Gaussian ones which are still widely used in financial engineering due to their convenience. Notice also the difference between asset classes: Though estimated correlations are $\rho = 0.34$ for the leftmost copulas, they have much dissimilar peculiarities.

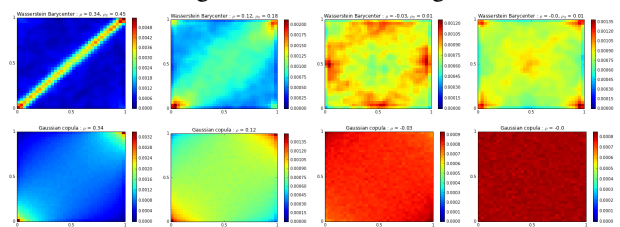
Stocks Centroids' main feature: More mass in the bottom-left corner, i.e. lower tail dependence. Stock prices tend to plummet together.



Credit default swaps Centroids' main feature: More mass in the top-right corner, i.e. upper tail dependence. Insurance cost against entities' default tends to soar in stressed market.



FX rates Centroids' main feature: Empirical copulas show that dependence between FX rates are various. For example, rates may exhibit either strong dependence or independence while being anti-correlated during extreme events.



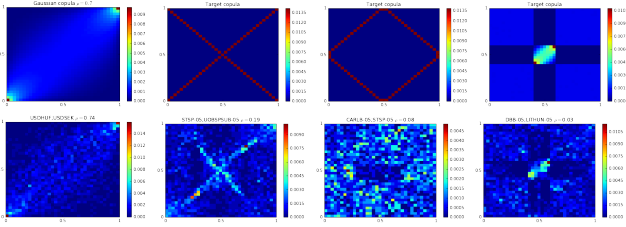


Figure 6: Target copulas (simulated or handcrafted) and their respective nearest copulas which answer questions A,B,C,D

Answering dependence queries

Inspired by the previous exploration results, we may want to answer such questions: (A) Which pair of assets having $\rho = 0.7$ correlation has the nearest copula to the Gaussian one? Though such questions can be answered by computing a likelihood for each pairs, our methodology stands out for dealing with non-parametric dependence patterns, and thus for questions such as: (B) Which pairs of assets are both positively and negatively correlated? (C) Which assets occur extreme variations while those of others are relatively small, and conversely? (D) Which pairs of assets are positively correlated for small variations but uncorrelated otherwise?

Considering a cross-asset dataset which comprises the SBF 120 components (index including the CAC 40 and 80 other highly capitalized French entities), the 500 most liquid CDS worldwide, and 80 FX rates, we display in Figure 6 the empirical copulas (alongside their respective targets) which best answer questions A,B,C,D.

Power of TFDC

In this experiment, we compare the empirical power of TFDC to well-known dependence coefficients such as Pearson linear correlation (cor), distance correlation (dCor) (Székely, Rizzo, and others 2009), maximal information coefficient (MIC) (Reshef et al. 2011), alternating conditional expectations (ACE) (Breiman and Friedman 1985), maximum mean discrepancy (MMD) (Gretton et al. 2012), copula maximum mean discrepancy (CMMMD) (Ghahramani, Póczos, and Schneider 2012), randomized dependence coefficient (RDC) (Lopez-Paz, Hennig, and Schölkopf 2013). Statistical power of a binary hypothesis test is the probability that the test correctly rejects the null hypothesis (H0) when the alternative hypothesis (H1) is true. In the case of dependence coefficients, we consider (H0): X and Y are independent; (H1): X and Y are dependent. Following the numerical experiment described in (Simon and Tibshirani 2014; Lopez-Paz, Hennig, and Schölkopf 2013), we estimate the power of the aforementioned dependence measures with simulated pairs of variables with different relationships (considered in (Reshef et al. 2011; Simon and Tibshirani 2014; Lopez-Paz, Hennig, and Schölkopf 2013)), but with varying levels of noise added. By design, TFDC aims at detecting the simulated dependence relationships. Thus, this dependence measure is expected to have a much higher power than coefficients such as MIC since, according to Simon and Tibshirani in (Simon and Tibshirani 2014), coef-

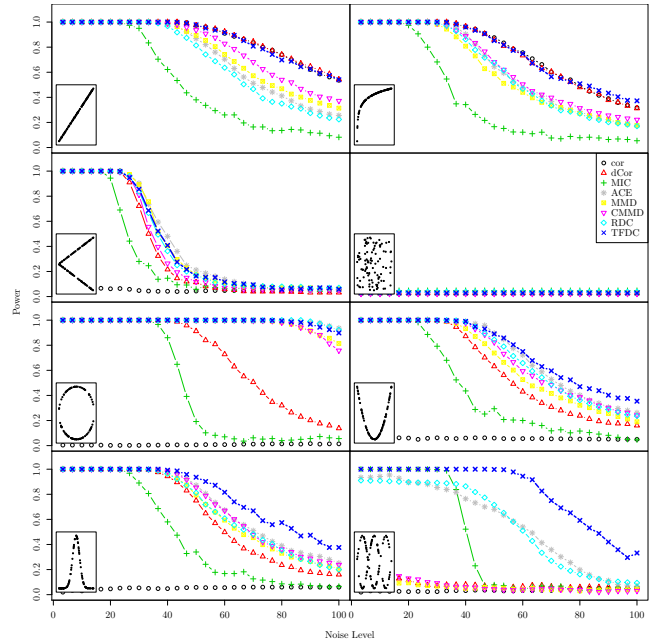


Figure 7: Power of several dependence coefficients as a function of the noise level in eight different scenarios. Insets show the noise-free form of each association pattern. The coefficient power was estimated via 500 simulations with sample size 500 each

ficients “which strive to have high power against all alternatives can have low power in many important situations.” TFDC only targets the specific important situations. Results are displayed in Figure 7.

Discussion

It is known by risk managers how dangerous it can be to rely solely on a correlation coefficient to measure dependence. That is why we have proposed a novel approach to explore, summarize and measure the pairwise correlations which exist between variables in a dataset. We have also pointed out through the UCI-datasets example that non-trivial dependence patterns can be easily found between the features variables. Using these patterns as *targets* when performing correlation-based feature selection may improve results. This idea still needs to be empirically verified. The experiments show the benefits of the proposed method: It allows to highlight the various dependence patterns that can be found between financial time series, which strongly depart from the Gaussian copula widely used in financial engineering. Though *answering dependence queries* as briefly outlined is still an art, we plan to develop a rich language so that a user can formulate complex questions about dependence, which will be automatically translated into copulas in order to let the methodology provide these questions accurate answers.

References

- Agueh, M., and Carlier, G. 2011. Barycenters in the Wasserstein space. *SIAM Journal on Mathematical Analysis* 43(2):904–924.
- Banerjee, A.; Merugu, S.; Dhillon, I. S.; and Ghosh, J. 2005. Clustering with Bregman divergences. *The Journal of Machine Learning Research* 6:1705–1749.
- Breiman, L., and Friedman, J. H. 1985. Estimating optimal transformations for multiple regression and correlation. *Journal of the American statistical Association* 80(391):580–598.
- Chang, Y.; Li, Y.; Ding, A.; and Dy, J. 2016. A robust-equitable copula dependence measure for feature selection. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, 84–92.
- Cuturi, M., and Doucet, A. 2014. Fast computation of Wasserstein barycenters. In *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21–26 June 2014*, 685–693.
- Cuturi, M. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, 2292–2300.
- Deheuvels, P. 1979. La fonction de dépendance empirique et ses propriétés. un test non paramétrique d'indépendance. *Acad. Roy. Belg. Bull. Cl. Sci.(5)* 65(6):274–292.
- Deheuvels, P. 1981. An asymptotic decomposition for multivariate distribution-free tests of independence. *Journal of Multivariate Analysis* 11(1):102–113.
- Dessein, A.; Papadakis, N.; and Rouas, J.-L. 2016. Regularized optimal transport and the rot mover's distance. *arXiv preprint arXiv:1610.06447*.
- Ding, A., and Li, Y. 2013. Copula correlation: An equitable dependence measure and extension of pearson's correlation. *arXiv preprint arXiv:1312.7214*.
- Durante, F.; Saminger-Platz, S.; and Sarkoci, P. 2009. Rectangular patchwork for bivariate copulas and tail dependence. *Communications in Statistics Theory and Methods* 38(15):2515–2527.
- Elidan, G. 2013. Copulas in machine learning. In *Copulae in mathematical and quantitative finance*. Springer. 39–60.
- Ghahramani, Z.; Póczos, B.; and Schneider, J. G. 2012. Copula-based kernel dependency measures. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, 775–782.
- Gretton, A.; Borgwardt, K. M.; Rasch, M. J.; Schölkopf, B.; and Smola, A. 2012. A kernel two-sample test. *Journal of Machine Learning Research* 13(Mar):723–773.
- Hall, M. A. 2000. Correlation-based feature selection for discrete and numeric class machine learning. In *Proceedings of the 17th International Conference on Machine Learning (ICML-00)*, 359–366.
- Kinney, J. B., and Atwal, G. S. 2014. Equitability, mutual information, and the maximal information coefficient. *Proceedings of the National Academy of Sciences* 111(9):3354–3359.
- Liebscher, E., et al. 2014. Copula-based dependence measures. *Dependence Modeling* 2(1):49–64.
- Lopez-Paz, D.; Hennig, P.; and Schölkopf, B. 2013. The randomized dependence coefficient. In *Advances in Neural Information Processing Systems*, 1–9.
- Marti, G.; Andler, S.; Nielsen, F.; and Donnat, P. 2016. Optimal transport vs. Fisher-Rao distance between copulas for clustering multivariate time series. In *2016 IEEE Statistical Signal Processing Workshop*.
- Monge, G. 1781. *Mémoire sur la théorie des déblais et des remblais*. De l'Imprimerie Royale.
- Muzellec, B.; Nock, R.; Patrini, G.; and Nielsen, F. 2016. Tsallis regularized optimal transport and ecological inference. *arXiv preprint arXiv:1609.04495*.
- Nelsen, R. B. 2013. *An introduction to copulas*, volume 139. Springer Science & Business Media.
- Reshef, D. N.; Reshef, Y. A.; Finucane, H. K.; Grossman, S. R.; McVean, G.; Turnbaugh, P. J.; Lander, E. S.; Mitzenmacher, M.; and Sabeti, P. C. 2011. Detecting novel associations in large data sets. *science* 334(6062):1518–1524.
- Reshef, D.; Reshef, Y.; Mitzenmacher, M.; and Sabeti, P. 2013. Equitability analysis of the maximal information coefficient, with comparisons. *arXiv preprint arXiv:1301.6314*.
- Rubner, Y.; Tomasi, C.; and Guibas, L. J. 2000. The earth mover's distance as a metric for image retrieval. *International journal of computer vision* 40(2):99–121.
- Sejdinovic, D.; Sriperumbudur, B.; Gretton, A.; Fukumizu, K.; et al. 2013. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *The Annals of Statistics* 41(5):2263–2291.
- Simon, N., and Tibshirani, R. 2014. Comment on "Detecting Novel Associations In Large Data Sets" by Reshef Et Al, Science Dec 16, 2011. *arXiv preprint arXiv:1401.7645*.
- Sklar, A. 1959. *Fonctions de répartition à n dimensions et leurs marges*. Université Paris 8.
- Székely, G. J.; Rizzo, M. L.; et al. 2009. Brownian distance covariance. *The annals of applied statistics* 3(4):1236–1265.
- Yu, L., and Liu, H. 2003. Feature selection for high-dimensional data: A fast correlation-based filter solution. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, 856–863.