

A proper scoring rule for minimum information copulas

Yici Chen* Tomonari Sei*†

April 11, 2022

Abstract

Multi-dimensional distributions whose marginal distributions are uniform are called copulas. Among them, the one that satisfies given constraints on expectation and is closest to the independent distribution in the sense of Kullback–Leibler divergence is called the minimum information copula. The density function of the minimum information copula contains a set of functions called the normalizing functions, which are often difficult to compute. Although a number of proper scoring rules for probability distributions having normalizing constants such as exponential families are proposed, these scores are not applicable to the minimum information copulas due to the normalizing functions. In this paper, we propose the conditional Kullback–Leibler score, which avoids computation of the normalizing functions. The main idea of its construction is to use pairs of observations. We show that the proposed score is strictly proper in the space of copula density functions and therefore the estimator derived from it has asymptotic consistency. Furthermore, the score is convex with respect to the parameters and can be easily optimized by the gradient methods.

Keywords: Copula, Homogeneity, Kullback–Leibler divergence, Multi-point locality, Normalizing function

1 Introduction

Multidimensional distributions whose marginal distributions are uniform are called copulas. Among them, the one that satisfies the expectation constraints and is closest to the independent distribution in the sense of Kullback–Leibler (KL) divergence is called the minimum information copula. Minimum information copulas are used in, e.g., financial models [5] and flood models [6]. It is known that the density function of a two-dimensional minimum information

*Graduate School of Information Science and Technology, The University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo, 113-8656, Japan.

†sei@mist.i.u-tokyo.ac.jp

copula is written in the following form [2]:

$$c_\theta(x, y) = \exp \left(\sum_{i=1}^k \theta_i h_i(x, y) + a_\theta(x) + b_\theta(y) \right),$$

where $h_i(x, y)$ ($1 \leq i \leq k$) are functions describing dependence between x and y , and $\theta = (\theta_i)$ is a parameter. The functions $a_\theta(x)$ and $b_\theta(y)$ are called normalizing functions that are determined by the marginal condition of copulas (see Section 2 for details). The normalizing functions play a similar role to the normalizing constant of exponential families. However, since the marginal condition involves a set of integral equations, it is generally more difficult to deal with the normalizing functions than the normalizing constants.

When performing estimation, a function that measures the goodness of the model, called the score function, is often used. Different scores have different properties, and by using scores with properties consistent with the intended use, estimation accuracy can be improved and the amount of computation can be reduced. For example, scores for estimation that are robust to noise have been proposed by [1] and [7] among others. The Hyvärinen score [9] has been proposed as a score without calculating the normalizing constant. A class of scoring rules with the same property are investigated by [13].

In this paper, we propose a scoring rule for minimum information copulas that can be calculated without the normalizing functions. The score, which we call the conditional Kullback–Leibler score, uses a conditional likelihood on pairs of observations. The score is shown to be strictly proper and therefore asymptotically consistent. Furthermore, the score is convex with respect to the parameters and can be easily optimized by the gradient methods.

The structure of this paper is as follows. Section 2 introduces copulas and minimum information copulas. Section 3 defines scores and explains their commonly used properties: propriety, locality, and homogeneity. In Section 4, we introduce general homogeneity and multi-point locality, which are key properties for dealing with normalizing functions. Then, we propose the conditional Kullback–Leibler score for minimum information copulas that satisfies generalized homogeneity and two-point locality. In Section 5, we confirm the asymptotic consistency through numerical experiments. Finally, future issues and prospects are discussed in Section 6.

For simplicity, we will only discuss the two-dimensional case, but we believe it can be extended to three or higher dimensions.

2 Minimum information copulas

2.1 Copulas

A copula is a joint distribution function C on the unit square with uniform marginals (e.g. [12]). In this paper, we only deal with absolutely continuous copulas, in which case the definition is stated by density functions as follows.

Definition 1 (Copula densities). A two-dimensional copula density is a function $c : [0, 1]^2 \rightarrow [0, \infty)$ that satisfies the following two properties:

$$\int_0^1 c(x, y) dy = 1, \quad x \in [0, 1] \quad (1)$$

and

$$\int_0^1 c(x, y) dx = 1, \quad y \in [0, 1]. \quad (2)$$

From Sklar's theorem, a joint distribution with arbitrary marginals can be constructed by a copula. Therefore, we can separate statistical modeling into the marginal part and the copula part.

Theorem 1 (Sklar's theorem). *Let h be a joint density function on \mathbb{R}^2 with marginal density functions f and g . Then there exists a copula density c such that for all $(x, y) \in \mathbb{R}^2$,*

$$h(x, y) = c(F(x), G(y))f(x)g(y),$$

where $F(x) = \int_{-\infty}^x f(\xi) d\xi$ and $G(y) = \int_{-\infty}^y g(\eta) d\eta$.

2.2 Minimum information copulas

Consider statistical modeling of a phenomenon with a multivariate data. Suppose that some constraints, such as mean and variance, are given as prior information. If there is little prior information about the distribution, the model that satisfies the constraints cannot be uniquely determined. In this case, which model should be adopted? One way to think of it is to adopt a neutral model that assumes as little as possible dependence between random variables that are not included in the prior information. In other words, we adopt the model that is closest to the distribution in which the random variables are independent.

The Kullback-Leibler(KL) divergence [11] is often used to measure the distance between distributions.

Definition 2 (Kullback–Leibler divergence). Let g_1 and g_2 be density functions. Then the Kullback–Leiber (KL) divergence is defined as follows:

$$D_{\text{KL}}(g_1||g_2) = \int g_1(x, y) \log \left(\frac{g_1(x, y)}{g_2(x, y)} \right) dx dy.$$

The copula that is closest to the independent copula in terms of the KL divergence is called the minimum information copula.

Definition 3 (Minimum information copulas). Let $h_1(x, y), \dots, h_k(x, y)$ be given functions and $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ be given numbers. Let $\pi(x, y) = 1$ be the independent copula density. Then, the copula density c that minimizes

$$D_{\text{KL}}(c||\pi) = \iint c(x, y) \log c(x, y) dx dy$$

subject to

$$\iint c(x, y)h_i(x, y)dxdy = \alpha_i \quad (1 \leq i \leq k)$$

is called the minimum information copula density.

The minimum information copula is characterized by the following theorem. See also [4] for details on the uniqueness and existence problem.

Theorem 2 ([2], Theorem 2 & Theorem 3). *The minimum information copula density is unique if it exists, and expressed in the following form:*

$$c(x, y) = \exp \left(\sum_{i=1}^k \theta_i h_i(x, y) + a(x) + b(y) \right)$$

with some θ_i , $a(x)$ and $b(y)$. The functions $a(x)$ and $b(y)$ are unique except for arbitrariness of the additive constants.

From the theorem, we can redefine the minimum information copula density by

$$c_\theta(x, y) = \exp \left(\sum_{i=1}^k \theta_i h_i(x, y) + a_\theta(x) + b_\theta(y) \right) \quad (3)$$

together with the marginal conditions (1) and (2). Here, the parameter of interest is $\theta = (\theta_i)_{i=1}^k$. The functions $a_\theta(x)$ and $b_\theta(y)$ are called the normalizing functions.

For identifiability of $c_\theta(x, y)$, we suppose that $h_1(x, y), \dots, h_k(x, y)$ are linearly independent modulo additive functions, which means that an identity

$$\sum_{i=1}^k \Theta_i h_i(x, y) + A(x) + B(y) = 0$$

for some Θ_i , $A(x)$ and $B(y)$ implies $\Theta_1 = \dots = \Theta_k = 0$.

3 Scoring rules

3.1 Scores

A score is a function used to measure the difference between a model and the true distribution. We only consider probability density functions on $[0, 1]^2$. Denote the set of probability density functions on $[0, 1]^2$ by \mathcal{P} .

Definition 4 (Scores). A score S is a real-valued function of $(x, y, q) \in [0, 1]^2 \times \mathcal{P}$. The expected value

$$S(p, q) = \iint S(x, y, q)p(x, y)dxdy$$

with respect to $p \in \mathcal{P}$ is called the expected score.

In this paper, the score function $S(x, y, q)$ and the expected score $S(p, q)$ use the same symbol S and are both called scores for convenience. One property that is necessary for a score to measure the difference between the true distribution and the model is called propriety.

Definition 5 (Propriety). A score S is said to be proper if $S(p, q) \geq S(p, p)$ for all density functions $p, q \in \mathcal{P}$.

In general, when we say score, we often refer to the proper score. A divergence can be defined from a proper score by

$$D(p||q) = S(p, q) - S(p, p),$$

which is like the distance between p and q . If p is fixed, minimizing the divergence and minimizing the score are equivalent.

Consider a statistical model $\{q_\theta\} \subset \mathcal{P}$ indexed by a parameter θ . If n data points (x_i, y_i) are observed and the empirical distribution is denoted as \hat{p} , the estimation can be operated as follows:

$$\hat{\theta} = \arg \min_{\theta} S(\hat{p}, q_\theta) = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n S(x_i, y_i, q_\theta). \quad (4)$$

There are many different types of scores, but one of the simplest is the local score.

Definition 6 (Locality [13]). A score S is said to be local (in strict sense) if there exists a function $s : [0, 1]^2 \times [0, \infty) \rightarrow \mathbb{R}$ such that

$$S(x, y, q) = s(x, y, q(x, y)).$$

Furthermore, for $l \geq 0$, a score S is said to be l -local if $S(x, y, q)$ is represented by at most the l -th derivatives of q at (x, y) .

Local scores are easy to calculate whenever $q(x, y)$ is explicitly expressed because the score can be obtained only from the information at that point, without integrating over the neighborhood or referring to other points.

Example 1 (KL score). The score

$$S(x, y, q) = -\log q(x, y)$$

is called the KL score because the divergence induced from it is the KL divergence. The KL score is 0-local and proper. In fact, it is known that the KL score is essentially the only score that is 0-local and proper [3].

3.2 Homogeneity

There are several computational advantages to using homogeneous scores for estimation.

Definition 7 (Homogeneity). A score S is said to be homogeneous if it satisfies $S(x, y, \lambda q) = S(x, y, q)$ for any constant $\lambda > 0$.

If a homogeneous score is used for estimation, computation of the normalizing constant is not necessary.

We have introduced the properties of scores: propriety, locality and homogeneity. Based on these definitions, we consider two examples of scores.

Example (Example 1 continued). The KL score is not homogeneous. Indeed, the KL score satisfies

$$S(x, y, \lambda q) = S(x, y, q) - \log \lambda.$$

Therefore, when the KL-score is used for estimation, computation of the normalizing constant is necessary.

Example 2 (Hyvärinen score [9]). A score

$$S(x, y, q) = \left(\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \log q(x, y) + \frac{1}{2} \left(\frac{\partial}{\partial x} \log q(x, y) \right)^2 + \frac{1}{2} \left(\frac{\partial}{\partial y} \log q(x, y) \right)^2$$

is called the Hyvärinen score. The Hyvärinen score is 2-local, homogeneous and proper. Therefore, the normalizing constant is not necessary for estimation. However, the Hyvärinen score is not useful for estimation of the minimum information copulas because it does not remove the normalizing functions.

4 The proposed score

4.1 General homogeneity

As mentioned in the last example, the normalizing functions $a_\theta(x)$ and $b_\theta(y)$ in the density function (3) do not vanish even if a homogeneous score is applied. For this reason, the following property is introduced.

Definition 8 (General homogeneity). A score S is said to be generally homogeneous if it satisfies that

$$\begin{aligned} S(x, y, \lambda_1 q) &= S(x, y, q), \\ S(x, y, \lambda_2 q) &= S(x, y, q) \end{aligned}$$

for any positive functions $\lambda_1(x)$ and $\lambda_2(y)$, where $\lambda_1 q$ and $\lambda_2 q$ are defined by $(\lambda_1 q)(x, y) = \lambda_1(x)q(x, y)$ and $(\lambda_2 q)(x, y) = \lambda_2(y)q(x, y)$, respectively.

Example 3. It is easy to see that a score

$$S(x, y, q) = -\frac{\partial^2}{\partial x \partial y} \log q(x, y)$$

is generally homogeneous and 2-local. However, the score is not proper. To see this, let $q(x, y)$ be the Gaussian density (over \mathbb{R}^2). Then $S(x, y, q)$ is a constant involving the correlation parameter and $S(p, q)$ takes any real value independent of p .

If S is generally homogeneous and $c_\theta(x, y)$ is the minimum information copula density in (3), we have

$$S(x, y, c_\theta(x, y)) = S(x, y, e^{\sum_i \theta_i h_i(x, y)}),$$

which does not require computation of the normalizing functions. Hence, our problem is reduced to find a generally homogeneous (l -)local proper score. However, after some trials based on symbolic computation in line with [13], the authors realized that such a score may not exist.

4.2 Multi-point scores and their locality

Instead of finding a generally homogeneous local proper score, we try to relax the required properties. General homogeneity and propriety are necessary for estimation of the minimum information copulas. Therefore, we reconsider locality. For this purpose, the concept of multi-point scores is introduced.

Definition 9 (Multi-point score). Let $m \geq 1$. An m -point score is a function

$$S(\mathbf{x}, \mathbf{y}, q) = S(x^1, y^1, \dots, x^m, y^m, q)$$

of $\mathbf{x} = (x^1, \dots, x^m) \in [0, 1]^m$, $\mathbf{y} = (y^1, \dots, y^m) \in [0, 1]^m$ and $q \in \mathcal{P}$. The expected score of $S(\mathbf{x}, \mathbf{y}, q)$ is defined as

$$S(p, q) = \int S(\mathbf{x}, \mathbf{y}, q) p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y},$$

where $p(\mathbf{x}, \mathbf{y}) = \prod_{i=1}^m p(x^i, y^i)$. In other words, the expectation is taken with respect to independent samples from p . A score is called a multi-point score if it is an m -point score for some m .

Definition of propriety and general homogeneity of the multi-point scores is straightforward. Locality is defined as follows.

Definition 10 (Multi-point locality). A multi-point score S is said to be local if there exists a function $s : [0, 1]^m \times [0, 1]^m \times [0, \infty)^{m \times m} \rightarrow \mathbb{R}$ such that

$$S(\mathbf{x}, \mathbf{y}, q) = s(\mathbf{x}, \mathbf{y}, \mathbf{q}),$$

where $\mathbf{q} = (q(x^i, y^j))_{i,j=1}^m$.

Locality defined in Section 3 meant that for a given model, the score at a single point is evaluated using only the information at that point. On the other hand, multi-point locality means that the score at m points is evaluated using all the information at the m^2 points $\{(x^i, y^j)\}_{i,j=1}^m$.

4.3 The conditional Kullback–Leibler score

In the following, we construct a generally homogeneous 2-point local proper score, which is applicable to estimation of minimum information copulas. Since it is known that the 1-point local proper score is essentially only the KL-score [3], which does not have general homogeneity, it is reasonable to construct 2-point local proper scores.

Definition 11 (The conditional Kullback–Leibler score). Define a 2-point local score by

$$S(x^1, y^1, x^2, y^2, q) = -\log\left(\frac{q^{11}q^{22}}{q^{11}q^{22} + q^{12}q^{21}}\right),$$

where $q^{ij} = q(x^i, y^j)$. We call it the conditional Kullback–Leibler score.

This score looks a little strange at first glance, but it can actually be seen as a kind of conditional KL score as follows.

Consider that there are two systems that are exactly the same and independent of each other. Data (x^1, y^1) are obtained from System 1, and data (x^2, y^2) are obtained from System 2. Suppose that, by mistake, we were able to record the numerical values of the data, but forgot the correspondence between x and y . Then, there are two possible cases: $(x^1, y^1), (x^2, y^2)$ or $(x^1, y^2), (x^2, y^1)$. Under the condition that only the numerical value of the data is known, the conditional probability that the data is the pair $(x^1, y^1), (x^2, y^2)$ is

$$\frac{q^{11}q^{22}}{q^{11}q^{22} + q^{12}q^{21}}.$$

Therefore, the score

$$S = -\log\left(\frac{q^{11}q^{22}}{q^{11}q^{22} + q^{12}q^{21}}\right)$$

will come naturally as the KL score for the conditional probability.

Now we state our main result.

Theorem 3. *The conditional KL score is generally homogeneous and proper.*

Proof. General homogeneity is straightforward. We prove the propriety as fol-

lows:

$$\begin{aligned}
& S(p, p) - S(p, q) \\
&= \int_0^1 \int_0^1 \int_0^1 \int_0^1 [S(x^1, y^1, x^2, y^2, p) - S(x^1, y^1, x^2, y^2, q)] p^{11} p^{22} dx^1 dy^1 dx^2 dy^2 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 [S(x^1, y^1, x^2, y^2, p) - S(x^1, y^1, x^2, y^2, q)] p^{11} p^{22} dx^1 dy^1 dx^2 dy^2 \\
&+ \frac{1}{2} \int_0^1 \cdots \int_0^1 [S(x^1, y^2, x^2, y^1, p) - S(x^1, y^2, x^2, y^1, q)] p^{12} p^{21} dx^1 dy^2 dx^2 dy^1 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\log \left(\frac{q^{11} q^{22}}{q^{11} q^{22} + q^{12} q^{21}} \right) - \log \left(\frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}} \right) \right] p^{11} p^{22} dx^1 dy^1 dx^2 dy^2 \\
&+ \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\log \left(\frac{q^{12} q^{21}}{q^{11} q^{22} + q^{12} q^{21}} \right) - \log \left(\frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}} \right) \right] p^{12} p^{21} dx^1 dy^2 dx^2 dy^1 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\log \frac{\frac{q^{11} q^{22}}{q^{11} q^{22} + q^{12} q^{21}}}{\frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}}} \right] \frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}} (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^1 dx^2 dy^2 \\
&+ \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\log \frac{\frac{q^{12} q^{21}}{q^{11} q^{22} + q^{12} q^{21}}}{\frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}}} \right] \frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}} (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^2 dx^2 dy^1.
\end{aligned}$$

Using the inequality $\log x \leq x - 1$, we obtain

$$\begin{aligned}
& S(p, p) - S(p, q) \\
&\leq \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\frac{\frac{q^{11} q^{22}}{q^{11} q^{22} + q^{12} q^{21}}}{\frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}}} - 1 \right] \frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}} (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^1 dx^2 dy^2 \\
&+ \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\frac{\frac{q^{12} q^{21}}{q^{11} q^{22} + q^{12} q^{21}}}{\frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}}} - 1 \right] \frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}} (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^2 dx^2 dy^1 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\frac{q^{11} q^{22}}{q^{11} q^{22} + q^{12} q^{21}} - \frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}} \right] (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^1 dx^2 dy^2 \\
&+ \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\frac{q^{12} q^{21}}{q^{11} q^{22} + q^{12} q^{21}} - \frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}} \right] (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^2 dx^2 dy^1 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 \left[\frac{q^{11} q^{22}}{q^{11} q^{22} + q^{12} q^{21}} - \frac{p^{11} p^{22}}{p^{11} p^{22} + p^{12} p^{21}} + \frac{q^{12} q^{21}}{q^{11} q^{22} + q^{12} q^{21}} - \frac{p^{12} p^{21}}{p^{11} p^{22} + p^{12} p^{21}} \right] \\
&\quad \times (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^1 dx^2 dy^2 \\
&= \frac{1}{2} \int_0^1 \cdots \int_0^1 [1 - 1] (p^{11} p^{22} + p^{12} p^{21}) dx^1 dy^1 dx^2 dy^2 \\
&= 0.
\end{aligned}$$

Therefore,

$$S(p, p) \leq S(p, q).$$

As a result, the score S is proper. Moreover, the equality condition is

$$\frac{q^{11}q^{22}}{q^{11}q^{22} + q^{12}q^{21}} = \frac{p^{11}p^{22}}{p^{11}p^{22} + p^{12}p^{21}} \quad (5)$$

for all (x^1, y^1, x^2, y^2) . \square

From the above, we construct a generally homogeneous 2-point local proper score for the general density function. In fact, if the target is restricted to minimum information copulas, this score has even stronger properties.

Definition 12 (Strict propriety). Let $\mathcal{M} \subset \mathcal{P}$ be a given class of probability density functions. A proper score S is said to be strictly proper relative to \mathcal{M} if the equality $S(p, q) = S(p, p)$ for $p, q \in \mathcal{M}$ implies $p = q$.

Theorem 4. *Let \mathcal{M} be a minimum information copula model. Then the conditional KL score is strictly proper relative to \mathcal{M} .*

Proof. Let $p, q \in \mathcal{M}$ and suppose that $S(p, q) = S(p, p)$. Then we have (5) in the proof of Theorem 3, which is equivalent to

$$\frac{q^{11}q^{22}}{q^{12}q^{21}} = \frac{p^{11}p^{22}}{p^{12}p^{21}},$$

that is,

$$\frac{q(x^1, y^1)q(x^2, y^2)}{q(x^1, y^2)q(x^2, y^1)} = \frac{p(x^1, y^1)p(x^2, y^2)}{p(x^1, y^2)p(x^2, y^1)}.$$

By fixing (x^1, y^1) to an arbitrary point, we obtain a relation

$$q(x^2, y^2) = p(x^2, y^2) \exp(a(x^2) + b(y^2)),$$

where $a(x^2) = \log(p^{11}q^{21}) - \log(q^{11}p^{21})$ and $b(y^2) = \log q^{12} - \log p^{12}$. Since p and q are assumed to be the minimum information copula densities, Theorem 2 implies

$$p(x^2, y^2) = q(x^2, y^2).$$

Therefore, the score is strictly proper relative to \mathcal{M} . \square

4.4 Properties of the estimator

We define an estimator based on the proposed score and briefly describe its properties.

For the conditional KL score, we first separate the given data randomly into $N = \lfloor n/2 \rfloor$ groups as

$$\{(x_i^1, y_i^1, x_i^2, y_i^2)\}_{i=1}^N.$$

Then, based on the empirical score

$$\hat{S}(\theta) = \frac{1}{N} \sum_{i=1}^N S(x_i^1, y_i^1, x_i^2, y_i^2, q),$$

the estimator is defined by

$$\hat{\theta} = \arg \min_{\theta} \hat{S}(\theta)$$

Recall the following theorem on consistency and asymptotic normality of estimators based on strictly proper scores. We omit regularity conditions to make the ideas clearer.

Theorem 5 ([8] and Theorem 5.23 of [15]). *Let θ_0 be the true parameter and suppose that $(x_1, y_1), \dots, (x_n, y_n)$ are independent and identically distributed. Let S be a strictly proper (1-point) score. Then, the estimator $\hat{\theta}$ defined by (4) converges almost surely to θ_0 as $n \rightarrow \infty$. Furthermore, under regularity conditions, the asymptotic normality holds:*

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N(0, J^{-1} V J^{-1}),$$

where $J = E[\nabla_{\theta} \nabla_{\theta}^{\top} S]$ and $V = E[(\nabla_{\theta} S)(\nabla_{\theta} S)^{\top}]$.

Since the conditional KL score is strictly proper as proved in Theorem 4, we obtain the following corollary.

Corollary 1. The estimator based on the conditional KL score is consistent and asymptotically normal.

Next we point out that the estimation is a convex optimization problem. Here, we use the vector notation $\boldsymbol{\theta} = (\theta_i)_{i=1}^k$ for convenience.

Theorem 6. *Consider a minimum information copula model*

$$q(x, y; \boldsymbol{\theta}) = \exp(\boldsymbol{\theta}^T \mathbf{h}(x, y) + a(x) + b(y)).$$

Suppose that $\mathbf{H}_1, \dots, \mathbf{H}_N \in \mathbb{R}^k$ are linearly independent, where

$$\mathbf{H}_i = \mathbf{h}_i^{12} + \mathbf{h}_i^{21} - \mathbf{h}_i^{22} - \mathbf{h}_i^{11}, \quad \mathbf{h}_i^{\alpha\beta} = \mathbf{h}(x_i^{\alpha}, y_i^{\beta}).$$

Then, the empirical score $\hat{S}(\boldsymbol{\theta})$ based on the conditional KL score is strictly convex with respect to $\boldsymbol{\theta}$.

Proof. It is easy to see

$$\begin{aligned}\hat{S} &= \frac{1}{N} \sum_{i=1}^N \left\{ -\log \left(\frac{q_i^{11} q_i^{22}}{q_i^{11} q_i^{22} + q_i^{12} q_i^{21}} \right) \right\} \\ &= \frac{1}{N} \sum_{i=1}^N \log \left(e^{\boldsymbol{\theta}^T \mathbf{H}_i} + 1 \right).\end{aligned}$$

Then, the gradient vector is

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} \hat{S} &= \frac{1}{N} \sum_{i=1}^N \frac{e^{\boldsymbol{\theta}^T \mathbf{H}_i}}{e^{\boldsymbol{\theta}^T \mathbf{H}_i} + 1} \mathbf{H}_i \\ &= \frac{1}{N} \sum_{i=1}^N \frac{\mathbf{H}_i}{e^{-\boldsymbol{\theta}^T \mathbf{H}_i} + 1}.\end{aligned}\tag{6}$$

The Hessian matrix

$$\frac{1}{N} \sum_{i=1}^N \frac{e^{-\boldsymbol{\theta}^T \mathbf{H}_i}}{(e^{-\boldsymbol{\theta}^T \mathbf{H}_i} + 1)^2} \mathbf{H}_i \mathbf{H}_i^T,$$

is positive definite under the assumption. Therefore, the score is strictly convex with respect to $\boldsymbol{\theta}$. \square

The estimation can be operated easily with the gradient method.

5 Numerical experiments

5.1 Experiments of Gaussian copulas

5.1.1 Setting

The normalizing function of a minimal information copula cannot be obtained in general. In this subsection, we consider the Gaussian copula, which is one of the few examples of minimal information copulas for which the normalization function can be obtained in a closed form.

Another reason for experimenting with a Gaussian copula is that the data can be easily generated by a variable transformation of the 2-dimensional normal distribution, and the maximum likelihood estimation (MLE) can be compared with the proposed method. In Subsection 5.2, we conduct numerical experiments for general minimum information copulas for which normalizing functions cannot be obtained.

The Gaussian copula is a multidimensional normal distribution with its marginal distributions converted to uniform distributions. The density function of the 2-dimensional Gaussian copula is

$$\frac{1}{2\pi(1-\rho^2)^{\frac{1}{2}}\phi(\xi)\phi(\eta)} \exp\left(-\frac{1}{2(1-\rho^2)}(\xi^2 - 2\rho\xi\eta + \eta^2)\right) \Big|_{\xi=\Phi^{-1}(x), \eta=\Phi^{-1}(y)}$$

for $(x, y) \in (0, 1)^2$, where ϕ and Φ are the density function and cumulative distribution function of the 1-dimensional standard normal variables ξ and η . The parameter ρ is the correlation coefficient of the Gaussian variables ξ and η .

As pointed out by [10], the Gaussian copula is in fact a minimum information copula with

$$\theta = \frac{\rho}{1 - \rho^2}$$

and

$$h(x, y) = \Phi^{-1}(x)\Phi^{-1}(y).$$

The normalizing functions are

$$\exp\{a(x)\} = \frac{1}{\sqrt{2\pi}(1 - \rho^2)^{\frac{1}{4}}\phi(\xi)} \exp\left(-\frac{\xi^2}{2(1 - \rho^2)}\right) \Big|_{\xi=\Phi^{-1}(x)}$$

and $\exp\{b(y)\} = \exp\{a(y)\}$. The parameters θ and ρ have one-to-one correspondence as

$$\rho = \frac{2\theta}{1 + \sqrt{1 + 4\theta^2}}.$$

Thus, the density function of the Gaussian copula can be expressed relatively simply. In the following, the procedure of the numerical experiment is explained in detail.

Setting. 1. First, we generate $2N$ data of 2-dimensional normal distribution with mean vector $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and covariance matrix $\begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$. Consider the first N data $\begin{bmatrix} \xi_i^1 \\ \eta_i^1 \end{bmatrix}$ as obtained from system 1 and the other N data $\begin{bmatrix} \xi_i^2 \\ \eta_i^2 \end{bmatrix}$ as obtained from system 2.

2. Take the variable transformation

$$\begin{bmatrix} x_i^j \\ y_i^j \end{bmatrix} = \begin{bmatrix} \Phi(\xi_i^j) \\ \Phi(\eta_i^j) \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\xi_i^j}{\sqrt{2}}\right)\right) \\ \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\eta_i^j}{\sqrt{2}}\right)\right) \end{bmatrix} \quad (7)$$

($i = 1, \dots, N; j = 1, 2$) in order to obtain $2N$ sample of the Gaussian copula, where $\operatorname{erf}(x) = (2/\sqrt{\pi}) \int_0^x e^{-t^2} dt$ is the error function.

3. Set the initial value $\theta = \theta_0$ and iterative step $d\theta$ properly.

4. Calculate the gradient of the empirical score

$$\hat{S} := \frac{1}{N} \sum_{i=1}^N S(x_i^1, y_i^1, x_i^2, y_i^2, q) = -\frac{1}{N} \sum_{i=1}^N \log \left(\frac{q_i^{11} q_i^{22}}{q_i^{11} q_i^{22} + q_i^{12} q_i^{21}} \right).$$

Substitute the sample $(x_i^j, y_i^j)(i = 1, \dots, N; j = 1, 2)$ for (6):

$$\frac{d}{d\theta} \hat{S} = \frac{1}{N} \sum_{i=1}^N \frac{H_i}{e^{-\theta H_i} + 1}.$$

5. If $|d\hat{S}/d\theta|$ is sufficiently small, output θ as the estimator $\hat{\theta}$. Otherwise, $\theta \leftarrow \theta - (d\hat{S}/d\theta)d\theta$ and go to Step 4.

The experiments of the MLE use the empirical score of the KL-score

$$\hat{S}_{\text{KL}} = -\frac{1}{N} \sum_{i=1}^N \log(q_i^{11} q_i^{22})$$

instead of the conditional KL score in Step 4. The other steps are the same.

5.1.2 Results

We describe the results of numerical calculations according to the experimental procedure in Subsection 5.1.1.

First, 4000 sample points of 2-dimensional normal distribution with mean vector $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and covariance matrix $\begin{bmatrix} 1 & \rho = 0.7 \\ \rho = 0.7 & 1 \end{bmatrix}$ were obtained. Then, the variable transformation (7) was operated in order to get the sample of Gaussian copula with parameter $\theta = \frac{\rho}{1-\rho^2} = 1.372549$. The larger the number of data, the more the estimation error converges to 0. To see the degree of convergence, the estimation error was calculated by changing the number of data used. Using $N = 40, 50, \dots, 1990, 2000$ pairs of data, we substituted the data into the proposal score and obtained the optimal solution $\hat{\theta}$ as the estimator. The estimation error was calculated as the absolute value of the difference from the true parameter $\theta = 1.372549$. In addition, maximum likelihood estimation was performed using the same data with KL-scores. The above experiment was repeated up to 100 times and the estimation errors were averaged.

The estimation results are shown in Figure 1. The results are also plotted on the logarithmic graph of both axes and linearly fitted in Figure 2. Linear fitting is the least-squares fitting of a, b of $y = x^a \exp(b)$ to the data. The red line is the estimation error of the maximum likelihood estimation (KL-score) and the green is the estimation error of the proposal score (CKL-score). The blue line is the estimation error fitting of the maximum likelihood estimation (KL-score), where the coefficients of the fitting are

$$\begin{aligned} a &= -0.517919 \\ b &= 0.445423. \end{aligned}$$

The purple line is the estimation error fitting of the proposed score (CKL-score), where

$$\begin{aligned} a &= -0.49438 \\ b &= 0.97278. \end{aligned}$$

The speed of convergence of the error is roughly $\frac{1}{\sqrt{N}}$ for the number N of data since the blue and purple a in the fitting are close to -0.5 . This result is consistent with Corollary 1. The maximum likelihood estimation has better results with respect to error convergence than the proposed method, as expected. However, when estimating parameters with minimum information copulas, the maximum likelihood estimation is only possible with the Gaussian copulas used in this experiment. In most other cases, the maximum likelihood estimation cannot be performed because the normalizing functions cannot be obtained. In such cases, the greatest strength of the proposed score is that it can estimate with the same accuracy.

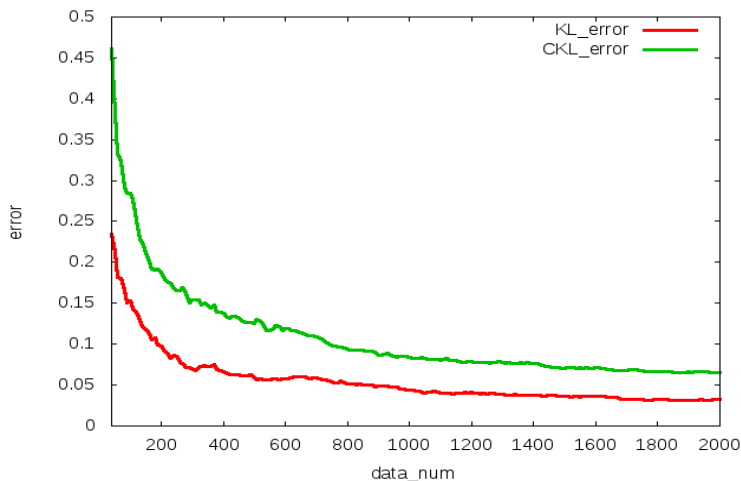


Figure 1: The estimation result of Gaussian copulas. The mean absolute error (vertical axis) with respect to the number of data (horizontal axis) is shown. The red line is the KL score (MLE) and green line is the CKL score (proposed)

5.2 Experiments of general minimum information copulas

5.2.1 Setting

Section 5.1 described the numerical experimental results of the Gaussian copulas. In this section, we present numerical experiments on general minimum information copulas. Since the normalizing functions of the minimum information copula are not generally obtainable, exact sampling is difficult. Thus, the generation of data itself is problematic before parameter estimation. In this paper, we generate data using the approximate sampling method proposed by [14]. The sampling procedure is described below.

Setting (The approximate sampling method). 1. First, we decide the function $h(x, y)$ and parameter $\theta \in \mathbb{R}$. Once these are determined, the min-

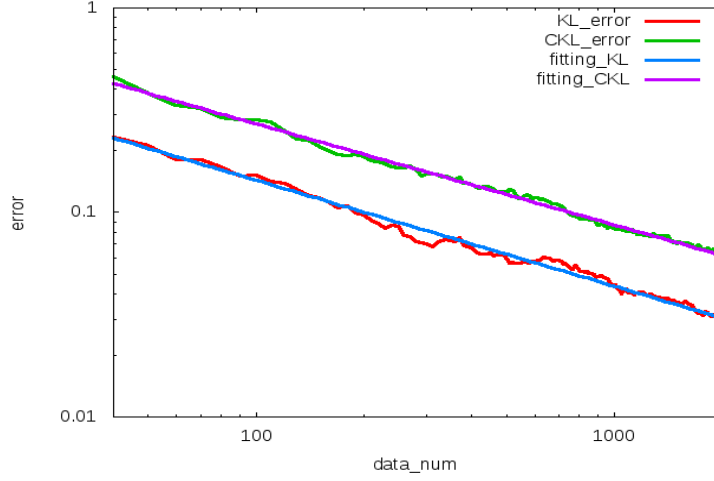


Figure 2: The estimation result of Gaussian copulas. The mean absolute error (vertical axis) with respect to the number of data (horizontal axis) is shown in logarithmic scale. The red line is the KL score (MLE), green line is the CKL score (proposed), and the blue and purple lines are the linear fitting of the two lines.

imum information copula is uniquely determined, although it cannot be written explicitly.

2. Then, generate a sufficiently large number of independent and identically distributed (i.i.d.) data $\begin{bmatrix} x_i \\ y_i \end{bmatrix}$ ($i = 1, \dots, 2N$) from the uniform distribution on $[0, 1]^2$.
3. Choose i, j that satisfy $1 \leq i < j \leq 2N$ randomly and flip them with probability

$$\begin{aligned} \rho &= \frac{q(x_i, y_j)q(x_j, y_i)}{q(x_i, y_i)q(x_j, y_j) + q(x_i, y_j)q(x_j, y_i)} \\ &= \frac{\exp\{\theta(h(x_i, y_j) + h(x_j, y_i))\}}{\exp\{\theta(h(x_i, y_i) + h(x_j, y_j))\} + \exp\{\theta(h(x_i, y_j) + h(x_j, y_i))\}} \end{aligned}$$

In other words, change the pair $\begin{bmatrix} x_i \\ y_i \end{bmatrix}$ and $\begin{bmatrix} x_j \\ y_j \end{bmatrix}$ to the pair $\begin{bmatrix} x_i \\ y_j \end{bmatrix}$ and $\begin{bmatrix} x_j \\ y_i \end{bmatrix}$.

4. Repeat step 3 enough times. Then, the $2N$ data are approximately an i.i.d. sample from the minimum information copula decided in step 1.
5. We consider the first N data $\begin{bmatrix} x_i^1 \\ y_i^1 \end{bmatrix}$ obtained from system 1 and the other N data $\begin{bmatrix} x_i^2 \\ y_i^2 \end{bmatrix}$ obtained from system 2.

Now we have an approximate sample from the given minimum information copulas and can perform parameter estimation using the proposed scores. Once the initial value θ_0 and the iterative step $d\theta$ are determined appropriately, the rest is done in the same way as in step 4 and 5 of Subsection 5.1.1. The above numerical experiments are repeated up to 100 times from sampling to estimation, and the average of the absolute errors is obtained.

5.2.2 Result

First, to roughly check the behavior of the approximate sampling, we conduct the experiment with Gaussian copulas again. The exact and approximate sampling were used, respectively, to obtain a sample of Gaussian copulas with the same parameters $\theta = \frac{\rho}{1-\rho^2} = 1.372549$ as in Subsection 5.1.2. Using $N = 40, 50, \dots, 1990, 2000$ pairs of data, we substituted the data into the proposed score and obtained the optimal solution $\hat{\theta}$ as the estimator. The estimation error was calculated as the absolute value of the difference from the true parameter $\theta = 1.372549$. In addition, maximum likelihood estimation was performed using the same data with KL-scores. The above experiment was repeated up to 100 times and the estimation errors were averaged.

The results of estimating data with exact and approximate sampling are shown in Figure 3, where both axes are in logarithmic scale. The red line shows the estimation error for the exact sampling and the maximum likelihood estimation (KL-score), and the green line is the estimation error for the exact sampling and the proposed score (CKL-score). The blue line is the estimation error of the approximate sampling and the maximum likelihood estimation (KL-score), and the purple line is the estimation error of the approximate sampling and proposed score (CKL-score).

The result of linear fitting for the four lines is summarized in Table 1. The figure and table confirm that the approximate sampling seems successfully generating data for the Gaussian copulas.

Table 1: The coefficients of linear fitting for the four lines on in Figure 3.

score	sampling	a	b
KL	exact	-0.517919	0.445423
CKL	exact	-0.49438	0.97278
KL	approximate	-0.492425	0.37061
CKL	approximate	-0.561812	1.52274

Since the rough behavior of approximate sampling has been confirmed, we will discuss the results of experiments on sampling and parameter estimation of general minimum information copulas, which was the original purpose of this paper.

We generated $N = 2000$ pairs of samples of minimum information copula with parameters $\theta = 5.0, 10.0$ and function $h(x, y) = xy, x^2y$. We used $N =$

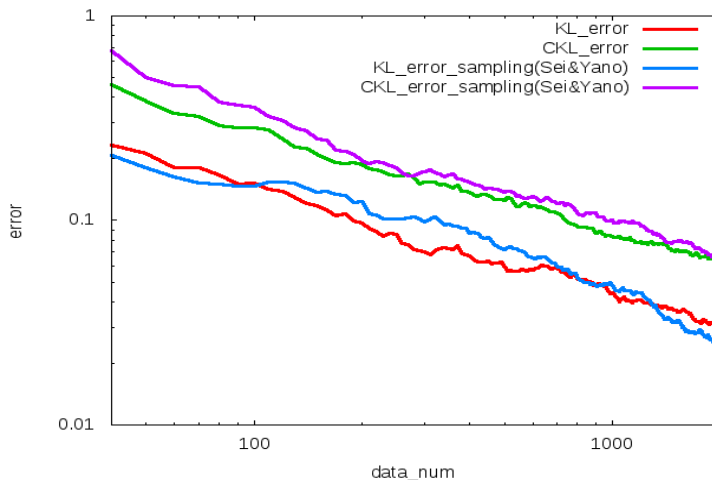


Figure 3: The estimation result of Gaussian copulas with exact and approximate sampling. The mean absolute error (vertical axis) with respect to the number of data (horizontal axis) is shown in logarithmic scale. The four lines mean the exact sampling + KL score (red), the exact sampling + CKL score (green), the approximate sampling + KL score (blue) and the approximate sampling + CKL score (purple).

20, 30, 40, \dots , 1990, 2000 pairs of them, substituted data into the proposed score, and estimated the optimal solution as $\hat{\theta}$. The estimation error was calculated as the absolute value of the difference from the true parameter θ . The above experiment was repeated up to 100 times and the average of the estimation error was taken.

The estimation results with the proposed score for the general minimum information copulas are shown in Figure 4, where both axes are in logarithmic scale. The red line is the estimation error for $\theta = 5.0, h(x, y) = xy$ and the green line is the estimation error for $\theta = 5.0, h(x, y) = x^2y$. The blue line is the estimation error for $\theta = 10.0, h(x, y) = xy$ and the purple is the estimation error for $\theta = 5.0, h(x, y) = x^2y$.

The result of linear fitting for the four lines is summarized in Table 2. The proposed method confirms that estimation is possible.

6 Conclusions

In this paper, we propose a generally homogeneous 2-point local strictly proper score for minimum information copulas. The greatest strength of this score is that it can be calculated without normalizing functions, which are difficult to compute for minimum information copulas. The estimator based on the score is asymptotically consistent, and numerical experiments have confirmed that

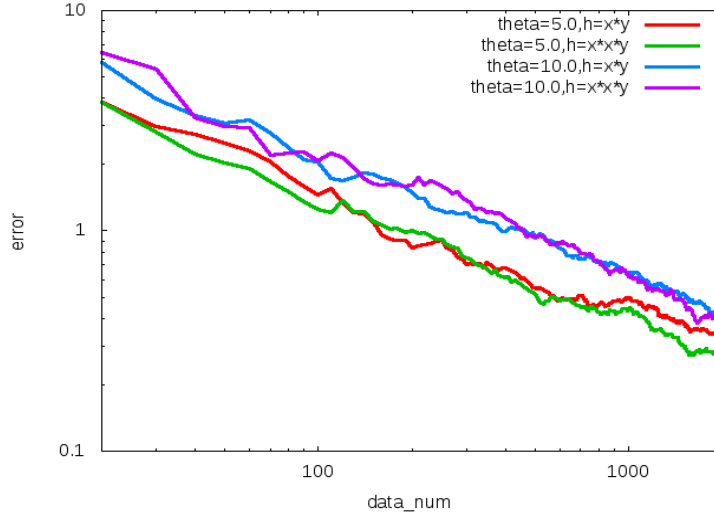


Figure 4: The estimation result of general minimum information copulas. The mean absolute error (vertical axis) with respect to the number of data (horizontal axis) is shown in logarithmic scale. The used parameters are $\theta = 5.0, h(x, y) = xy$ (red), $\theta = 5.0, h(x, y) = x^2y$ (green), $\theta = 10.0, h(x, y) = xy$ (blue) and $\theta = 10.0, h(x, y) = x^2y$ (purple).

Table 2: The coefficients of linear fitting for the four lines on in Figure 4.

θ	$h(x, y)$	a	b
5.0	xy	-0.563269	3.02252
5.0	x^2y	-0.572157	2.98715
10.0	xy	-0.548056	3.30813
10.0	x^2y	-0.581245	3.52886

the behavior is consistent with the theory. Future work includes theoretical computation of asymptotic variance.

In addition, in this paper, we considered the data

$$(x_1, y_1), \dots, (x_n, y_n)$$

as $N = \lfloor n/2 \rfloor$ pairs

$$(x_1^1, y_1^1, x_1^2, y_1^2), \dots, (x_N^1, y_N^1, x_N^2, y_N^2).$$

However, there are $n(n-1)/2$ pairs in the data: (x_i, y_i, x_j, y_j) for $1 \leq i < j \leq n$. In terms of extracting the full information of the data, it would be better to consider the empirical score of all combinations

$$\hat{S} = \frac{2}{n(n-1)} \sum_{1 \leq i < j \leq n} \log \left(\frac{q_{ii}q_{jj}}{q_{ii}q_{jj} + q_{ij}q_{ji}} \right), \quad q_{ij} := q(x_i, y_j).$$

Comparison of the accuracy and computational speed of these two approaches are quite interesting. The theory of U-statistics (e.g. Chapter 12 of [15]) may be useful for analysis of the modified score.

Furthermore, other than the proposed score, scores with general homogeneity, propriety, and multi-point locality should also be discussed. In this paper, the score was based on the form of the KL score, but it is not yet known whether a score with similar properties can be created by mimicking the form of the Hyvärinen score, for example. It is also possible that scores with similar properties can be created using a form that does not take logarithms. Thus, the authors believe that there is still a way to create such a score. By knowing the entire set of scores with these properties, we can see if there are scores among them that are even more accurate or that satisfy the properties we want to add.

Acknowledgments

We would like to thank Keisuke Yano for helpful comments.

References

- [1] A. Basu, I. R. Harris, N. L. Hjort, and M. C. Jones. Robust and efficient estimation by minimising a density power divergence. *Biometrika*, 85(3):549–559, 09 1998.
- [2] T. Bedford and K. Wilson. On the construction of minimum information bivariate copula families. *Annals of the Institute of Statistical Mathematics*, 66, 08 2014.
- [3] J. M. Bernardo. Expected information as expected utility. *The Annals of Statistics*, 7(3):686 – 690, 1979.

- [4] J. M. Borwein, A. S. Lewis, and R. D. Nussbaum. Entropy minimization, DAD problems, and doubly stochastic kernels. *J. Funct. Anal.*, 123:264–307, 1994.
- [5] O. Chatrabgoun, A. Hosseinian-Far, V. Chang, N. G. Stocks, and A. Daneshkhah. Approximating non-Gaussian Bayesian networks using minimum information vine model with applications in financial modelling. *Journal of Computational Science*, 24:266–276, 2018.
- [6] A. Daneshkhah, R. Remesan, O. Chatrabgoun, and I. P. Holman. Probabilistic modeling of flood characterizations with parametric and minimum information pair-copula model. *Journal of Hydrology*, 540:469–487, 2016.
- [7] H. Fujisawa and S. Eguchi. Robust parameter estimation with a small bias against heavy contamination. *Journal of Multivariate Analysis*, 99(9):2053–2081, 2008.
- [8] T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.
- [9] A. Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(Apr):695–709, 2005.
- [10] M. J. Jansen. Maximum entropy distributions with prescribed marginals and normal score correlations. In *Distributions with given marginals and moment problems*, pages 87–92. Springer, 1997.
- [11] S. Kullback and R. A. Leibler. On information and sufficiency. *The Annals of Mathematical Statistics*, 22(1):79 – 86, 1951.
- [12] R. B. Nelsen. *An Introduction to Copulas*. Springer, New York, NY, USA, second edition, 2006.
- [13] M. Parry, A. P. Dawid, and S. Lauritzen. Proper local scoring rules. *Ann. Statist.*, 40(1):561–592, 02 2012.
- [14] T. Sei and K. Yano. Minimum information dependence modeling for arbitrary product spaces. in preparation, 2022.
- [15] A. W. V. d. Vaart. *Asymptotic Statistics*. Cambridge University Press, 2000.