# Sobolev convergence of empirical Bernstein copulas

Sundusit Saekow[1] , Santi Tasena*,[2],[3]

[1] *Graduate Degree Program in Mathematics, Faculty of Science, Chiang Mai University, Chiang Mai, 50200, Thailand.*
[2] *Excellence in Mathematics, CHE, Bangkok 10400, Thailand.*
[3] *Department of Mathematics, Faculty of Science, Chiang Mai University, Chiang Mai, 50200 Thailand.*

## Abstract

In this work, we prove that Bernstein estimator always converges to the true copula under Sobolev distances. The rate of convergences is provided in case the true copula has bounded second order derivatives. Simulation study has also been done for Clayton copulas. We then use this estimator to estimate measures of complete dependence for weather data. The result suggests a nonlinear relationship between the dust density in Chiang Mai, Thailand and the temperature and the humidity level.

**Mathematics Subject Classification (2010).** 62H12, 62H20

**Keywords.** Bernstein estimator, copula, empirical, Sobolev convergence

## 1. Introduction

Sklar's Theorem[15] states that any joint distribution function can be written as a composition of a copula and its marginal distribution functions. In this sense, copulas represent links between random variables. Copulas model are then use in several fields such as economics, finance, drought, hydrology, etc.

In practice, the true copula is unknown; it has to be estimated from data. This situation is similar to that of (joint) distribution functions. It is well-known that the empirical copula always converges to the true copula under the uniform metric which proved to be sufficient in most cases. This type of convergence, never the less, is not sufficient for the estimation of measures of (mutual) complete dependence.

The concept of complete dependence can be traced back to Rényi [11] and Lancaster [9]. Both of them propose axioms for measuring association between random variables. These axioms include (1) the measure should be defined for all continuous random variables with its range between zero and one, (2) the measure should be zero only when the random variables are independent from each other, and (3) the measure should be invariant of the marginal distributions of random variables, in other words, it should be independent on rescaling of random variables. Rényi additionally requires that the measure takes the extreme value one if and only if both random variables are completely dependence of one

---
*Corresponding Author.
Email addresses: sundusit_saekow@cmu.ac.th (S. Saekow), santi.tasena@cmu.ac.th (S. Tasena)

another, viz., one is a bijective function of another while Lancaster only require that one is a function of another and not necessary vice versa.

Similar to the classical Pearson correlation coefficient and the Spearman's rank correlation coefficient, measures of complete dependence can be used to detect relationship between random variables. The difference is that the Pearson correlation coefficient only reach the maximum value when one random variable is a linear function of another and the Spearman's rank correlation coefficient only reach the maximum value when one random variable is a monotone function of another. Thus, these two measures might miss nonlinear relationship between these random variables. For example, it might be possible that both of these coefficients are roughly zero but the measures of complete dependence yield positive value.

Several decades later, measures of complete dependence has been proposed based on copulas[3, 14, 18]. Even though these measures are functions of copulas, they are not continuous with respect to the uniform convergence. Thus, the fact that empirical copulas converge to the true copula uniformly will not imply the convergences of these measures. Moreover, these measures are defined in term of derivative of copulas which implies a modified version of empirical copulas is needed since their partial derivatives do not exist.

Recently, Janssen et al.[6] consider the Bernstein estimator instead of the empirical copulas and they proved that this estimator also converges uniformly to the true copula. Also, this estimator is differentiable. Thus, it is natural to ask whether this estimator can be used to estimate measures of complete dependence. After all, the Bernstein copulas converge to the true copula under the Sobolev distance [16, Theorem 3]. This result does not directly implies the Bernstein empirical copulas converge to the true copula under the Sobolev distance, nevertheless. This is because the latter is based on the empirical copula which is varied with statistical samples. This fact also implies that the latter can be computed and be used as an estimator while the former is unknown because it is based on the true copula. Note that an estimator based on kernel method has also been defined for the measure in [3] for copulas which are thrice continuously differentiable.

In this work, we show that the Bernstein estimator converges to the true copula under Sobolev distance regardless of the differentability of the true copulas. The rate of convergences can also be computed as long as that of Bernstein copulas is known. We also compute the rate of convergence in case the true copula has bounded second order derivatives. Our result weaken the assumption of the kernel estimator in [3] but at the same time also weaken the conclusion. Here we only prove the law of large number while in [3], the central limit theorem is proven. We also provide numerical study of our estimator. The choice of copulas is chosen to match that of [3] for comparison.

Computation of measures of dependence for weather data is also shown. The dust density have significantly increase in Northern Thailand in recent years. It is a challenge to understand the dust behavior and to control it. Several censors have been installed to collect data. The weather data used in this work is taken from the Climate Change Data Center of Chiang Mai University[†]. The 966 sets of hourly average open data from 6 censors located within the Mueang District, Chiang Mai between September 14th, 2018 and September 25th, 2018 are used to compute the value of measures of complete dependence. The location of these censors are given in Table 1.

The source variable considered is either the temperature or the humidity level and the target variable is the dust density ($ug/m^3$) in the air. The Pearson correlation and the Spearman's rank correlation are also computed for comparison (see Table 11). We confirm that the dust density depends on the humidity level and the temperature. Their relationship is nonlinear, however, due to the fact that both Pearson correlation and the Spearman's rank correlation are small. Unfortunately, the current data is collected during

---

[†]website: www.cmuccdc.org

**Table 1.** Censor locations

| Location | Latitude | Longitude |
|---|---|---|
| Pa Dad Subdistrict | 18.745484 | 98.980442 |
| CMU, Mae Hia | 18.761371 | 98.931855 |
| Chiang Mai Chamber | 18.791512 | 99.018145 |
| Ban Thammapakorn | 18.782758 | 98.993233 |
| Night Bazaar | 18.7854765 | 98.9996602 |
| CMRU | 18.804701 | 98.986802 |

low dust density period. These values will have to be recomputed again when the data is ready in the future. An estimator in the case that the source variable is a random vector, such as temperature and humidity level, also has to be constructed.

The organization of this work is as follows. The next section provides basic concepts and terminologies needed for this work. Section 3 provides the proof of convergences of Bernstein estimator. Section 4 provides numerical results of this estimator.

## 2. Basic concepts and terminologies

Henceforth, let $\mathbb{I}$ denote the unit interval and $\mathbb{R}$ denote the set of real numbers. A copula is simply a (bivariate) joint distribution function with uniform marginals restricted to $\mathbb{I}^2$. A subcopula is then a restriction of a copula on some closed subset $\mathbb{A} \times \mathbb{B}$ of $\mathbb{I}^2$ containing $\{0,1\}^2$. Analytically, a subcopula is a function $S : \mathbb{A} \times \mathbb{B} \to \mathbb{I}$ such that

(1) $S$ is grounded, viz., $S(0,v) = 0 = S(u,0)$ for all $u \in \mathbb{A}$ and $v \in \mathbb{B}$,
(2) $S(1,v) = v$ and $S(u,1) = u$ for all $u \in \mathbb{A}$ and $v \in \mathbb{B}$, and
(3) $S(u_2,v_2) - S(u_2,v_1) - S(u_1,v_2) + S(u_1,v_1) \geq 0$ whenever $(u_i,v_i) \in \mathbb{A} \times \mathbb{B}$, $u_1 \leq u_2$, and $v_1 \leq v_2$.

A copula is then a subcopula with domain $\mathbb{I}^2$. Also, any copula is differentiable a.e. and its derivative is between zero and one. The usage of copulas in modeling dependence structure is due to the following Sklar's Theorem.

**Sklar's Theorem.** *For any joint distribution $H : \mathbb{R}^2 \to \mathbb{I}$ with marginal distributions $F : \mathbb{R} \to \mathbb{I}$ and $G : \mathbb{R} \to \mathbb{I}$, there is a copula $C : \mathbb{I}^2 \to \mathbb{I}$ such that*

$$H(x,y) = C(F(x), G(y))$$

*for all $x, y \in \mathbb{R}$. Moreover, this copula $C$ is unique if both $F$ and $G$ are continuous.*

According to Sklar's Theorem, copulas capture relationships between random variables while their marginal distributions capture individual behaviors. Thus, a measure of dependence should be defined in term of copulas if it is to be invariant under rescaling of random variables. This leads to the construction of several copula-based measures of dependence [1–5,7,8,10,12–14,17,19,20]. Most of these measures are continuous with respected to the uniform convergence excepted for measures of complete dependence defined in [3,14,18].

Truschnig [18] defines a measure of complete dependence $\zeta_1 (Y|X) = \zeta_1 (C_{X,Y})$ based on the Sobolev $L^1$-distance via

$$\zeta_1 (C_{X,Y}) = 3 \int_{\mathbb{I}} \int_{\mathbb{I}} |\partial_u C_{X,Y}(u,v) - v| \, du dv$$

while Siburg and Stoimenov[14] define a measure of complete dependence $\omega (X,Y) = \omega (C_{X,Y})$ based on the Sobolev $L^2$-distance via

$$\omega (C_{X,Y}) = \sqrt{3 \left( \int_{\mathbb{I}} \int_{\mathbb{I}} |\partial_u C_{X,Y}(u,v)|^2 \, du dv + \int_{\mathbb{I}} \int_{\mathbb{I}} |\partial_v C_{X,Y}(u,v)|^2 \, du dv \right) - 2}$$

where $C_{X,Y}$ is the copula associated with the joint distribution of (continuous) random variables $X$ and $Y$ as in the Sklar's Theorem. Note that an asymmetric version of $\omega$ is also defined by [3] via

$$r\left(Y|X\right) = r\left(C_{X,Y}\right) = 6\int_{\mathbb{I}}\int_{\mathbb{I}}\left|\partial_u C_{X,Y}(u,v)\right|^2 dudv - 2.$$

It can be easily seen that

$$\omega\left(X,Y\right) = \sqrt{\frac{1}{2}\left(r\left(Y|X\right) + r\left(X|Y\right)\right)}.$$

These measures all have ranges between zero and one. Also, $\zeta_1\left(Y|X\right)$ and $r\left(Y|X\right)$ reach the extreme value one if and only if $Y$ is a function of $X$. Since the set of such copulas is dense in the space of copula under the uniform convergence, these measures can not be continuous under the uniform metric. Otherwise, these measures would take constant value one.

In practice, the true copula associated with random variables $X$ and $Y$ are unknown and has to be estimated from data. Let $(X_1, Y_1), \ldots, (X_n, Y_n)$ be an i.i.d. sample of $(X, Y)$. Recall that the empirical joint distribution function of $(X, Y)$ is defined by

$$H_n\left(x,y\right) = \frac{1}{n}\sum_{i=1}^{n} 1_{\{X_i \le x, Y_i \le y\}}$$

for all $x, y \in \mathbb{R}$. The empirical marginal distribution functions $F_n$ and $G_n$ of $X$ and $Y$ are defined similarly. The empirical copula is then defined, for all $u, v \in \mathbb{I}$, by

$$C_n(u,v) = H_n\left(F_n^-(u), G_n^-(v)\right) = \frac{1}{n}\sum_{i=1}^{n} 1_{\{F_n(X_i) \le u, G_n(Y_i) \le v\}}$$

where $F^-$ stands for the quantile function associated with the distribution function $F$. It has been proved that $C_n \to C_{X,Y}$ uniformly.

**Theorem 2.1.** *[6, Lemma 1] Let $C_n$ be the empirical copula of the copula $C = C_{X,Y}$ as defined above. Then*

$$d_\infty\left(C_n, C\right) = \sup_{u,v \in \mathbb{I}} |C_n(u,v) - C(u,v)| = O\left(\sqrt{\frac{\ln\ln n}{n}}\right) \quad a.s.$$

*as $n \to \infty$.*

Note that $C_n$ is not a (random) copula. It is not even a (random) distribution function. If $C_n$ is restricted on $\left\{0, \frac{1}{n}, \ldots, \frac{n-1}{n}, 1\right\}^2$, however, it is a subcopula. Janssen et al. [6] use the idea of Bernstein copulas to modify $C_n$ so that it becomes copula and prove the convergence theorem. Recall that the Bernstein copula $B_m(S)$ of a subcopula $S$ in which its domain contains $\left\{0, \frac{1}{m}, \ldots, \frac{m-1}{m}, 1\right\}^2$ is defined by

$$B_m(S)(u,v) = \sum_{i=0}^{m}\sum_{j=0}^{m} S\left(\frac{i}{m}, \frac{j}{m}\right) P_{i,m}(u) P_{j,m}(v)$$

for all $u, v \in \mathbb{I}$ where $P_{i,m}(u) = \binom{m}{i} u^i (1-u)^{m-i}$ is the mass function of the binomial distribution. Any Bernstein copula is always a copula and it is smooth on the interior of $\mathbb{I}^2$. Moreover, $B_m(C) \to C$ under the Sobolev distance for all copula $C$ [16]. In general, $B_m(A)$ can be defined for any function $A : \mathbb{I}^2 \to \mathbb{R}$ although $B_m(A)$ might not be a copula in this case. Still, we have

$$\partial_u B_m\left(A\right)(u,v) = m\sum_{i=0}^{m-1}\sum_{j=0}^{m}\left(A\left(\frac{i+1}{m}, \frac{j}{m}\right) - A\left(\frac{i}{m}, \frac{j}{m}\right)\right) P_{i,m-1}(u) P_{j,m}(v)$$

for all $u, v \in \mathbb{I}$ [16, Equation 3].

The empirical Bernstein copula is defined to be

$$C_{m,n} = B_m\left(C_n\right)$$

where $C_n$ is the empirical copula. Note that $C_{m,n}$ is a (random) copula when $m$ divides $n$. Also, $0 \leq \partial_u C_{m,n}(u, v) \leq 2$ when $m \leq n$ regardless of whether $m$ divides $n$. This follows from the fact that

$$
\begin{aligned}
0 &\leq C_n\left(\frac{i+1}{m}, \frac{j}{m}\right) - C_n\left(\frac{i}{m}, \frac{j}{m}\right) \\
&= C_n\left(\frac{1}{n}\left\lceil\frac{n(i+1)}{m}\right\rceil, \frac{j}{m}\right) - C_n\left(\frac{1}{n}\left\lceil\frac{ni}{m}\right\rceil, \frac{j}{m}\right) \\
&\leq \frac{1}{n}\left\lceil\frac{n(i+1)}{m}\right\rceil - \frac{1}{n}\left\lceil\frac{ni}{m}\right\rceil \\
&\leq \frac{1}{n}\left(\frac{n(i+1)}{m} + 1 - \frac{ni}{m}\right) \\
&= \frac{m+n}{mn}.
\end{aligned}
$$

The following result provides a rate in which the Bernstein empirical copula $C_{m,n}$ uniformly converges to the true copula $C$.

**Theorem 2.2.** *[6, Theorem 1] If $m$ is a function of $n$ for which $m = m(n) \to \infty$ and $\frac{n}{m \ln \ln n} \to c < \infty$, then*

$$d_\infty\left(C_{m,n}, C\right) = O\left(\sqrt{\frac{\ln \ln n}{n}}\right) \quad a.s.$$

*as $n \to \infty$.*

Henceforth, denote $\|\cdot\|_p$ the (modified) Sobolev $L^p$-norm, that is,

$$\|A\|_p = \left(\int_\mathbb{I}\int_\mathbb{I} |\partial_u A|^p \, du dv\right)^{1/p}$$

for all differentiable function $A : \mathbb{I}^2 \to \mathbb{R}$. Note that $\|\cdot\|_p$ is an actual norm on the vector space spanned by copulas since all copulas are grounded. Moreover,

$$\|A - B\|_p^p \leq \|A - B\|_1 \leq \|A - B\|_p$$

for all copulas $A$ and $B$ since their partial derivatives lie in $\mathbb{I}$. Thus, all Sobolev $L^p$-norm induce the same topology. It is also known that the Sobolev norm is stronger than uniform norm but these two are equivalent since the set of shuffles of min is dense in the space of copulas under the uniform norm but it is nowhere dense under the Sobolev norm.

In the next section, we will show that $C_{m,n}$ also converges to the true copula $C$ under Sobolev norm. We will also provide the rate of convergence in case $C$ has bounded second order derivatives. The proofs rely on the following simple facts regarding the binomial distribution.

**Theorem 2.3.** *For any $i = 0, \ldots, m$ and any $0 < u < 1$,*

(1) $\partial_u P_{i,m}(u) = \frac{i - mu}{u(1-u)} P_{i,m}(u)$,

    (a) $\displaystyle\sum_{i=0}^{m} (i - mu)^2 P_{i,m}(u) = mu(1-u)$,

    (b) $\displaystyle\sum_{i=0}^{m} (i - mu)^4 P_{i,m}(u) = mu(1-u)\left((3m-6)u(1-u) + 1\right)$.

## 3. Sobolev convergence of Bernstein estimator

Using the same notations as in the previous section, we will prove that $\|C_{m,n} - C\|_p \to 0$ a.s.

**Lemma 3.1.** *If $m$ is a function of $n$ for which $m = m(n) \to \infty$ and $\frac{n}{m \ln \ln n} \to c < \infty$, then*

$$\mathbb{E}d_\infty\left(C_{m,n}, C\right) = O\left(\sqrt{\frac{\ln \ln n}{n}}\right)$$

*as $n \to \infty$.*

The proof of this lemma is actually adapted from [6, Lemma 3] with the help of the law of iterated logarithm.

**Proof.** Let $U_i = F(X_i)$ and $V_i = G(Y_i)$. Then $(U_i, V_i)$ are i.i.d. with the copula $C$ as their joint distribution function. Denote $\bar{H}_n$ the empirical distribution function of $C$ defined using $(U_i, V_i)$ as the sample, and let $\bar{F}_n$ and $\bar{G}_n$ be the marginals of $\bar{H}_n$. Then $\bar{F}_n^- = F F_n^-$ and $\bar{G}_n^- = G G_n^-$ which implies $C_n(u, v) = \bar{H}_n\left(\bar{F}_n^-(u), \bar{G}_n^-(v)\right)$ for all $u, v \in \mathbb{I}$. Thus,

$$\begin{aligned}
d_\infty\left(C_n, C\right) &= \sup_{u,v \in \mathbb{I}} \left| \bar{H}_n\left(\bar{F}_n^-(u), \bar{G}_n^-(v)\right) - C(u, v) \right| \\
&\leq \sup_{u,v \in \mathbb{I}} \left| \bar{H}_n\left(\bar{F}_n^-(u), \bar{G}_n^-(v)\right) - C(\bar{F}_n^-(u), \bar{G}_n^-(v)) \right| \\
&\quad + \sup_{u \in \mathbb{I}} \left| \bar{F}_n^-(u) - u \right| + \sup_{v \in \mathbb{I}} \left| \bar{G}_n^-(v) - v \right| \\
&\leq d_\infty\left(\bar{H}_n, C\right) + \sup_{u \in \mathbb{I}} \left| U_{(\lceil nu \rceil)} - \frac{1}{n}\lceil nu \rceil \right| + \sup_{u \in \mathbb{I}} \left| u - \frac{1}{n}\lceil nu \rceil \right| \\
&\quad + \sup_{v \in \mathbb{I}} \left| V_{(\lceil nv \rceil)} - \frac{1}{n}\lceil nv \rceil \right| + \sup_{u,v \in \mathbb{I}} \left| v - \frac{1}{n}\lceil nv \rceil \right| \\
&\leq d_\infty\left(\bar{H}_n, C\right) + \sup_{u \in \mathbb{I}} \left| U_{(\lceil nu \rceil)} - \bar{F}_n\left(U_{(\lceil nu \rceil)}\right) \right| \\
&\quad + \sup_{v \in \mathbb{I}} \left| V_{(\lceil nv \rceil)} - \bar{G}_n\left(V_{(\lceil nv \rceil)}\right) \right| + \frac{2}{n} \\
&\leq d_\infty\left(\bar{H}_n, C\right) + d_\infty\left(\bar{F}_n, Id\right) + d_\infty\left(\bar{G}_n, Id\right) + \frac{2}{n} \\
&\leq 3d_\infty\left(\bar{H}_n, C\right) + \frac{2}{n}
\end{aligned}$$

where $Id : \mathbb{I} \to \mathbb{I}$ is the identity function. By law of iterated logarithm,

$$\limsup_{n \to \infty} \frac{\sqrt{n}d_\infty\left(C_n, C\right)}{\sqrt{2 \ln \ln n}} \leq \frac{3}{2} \text{ a.s.}$$

which implies $\limsup_{n \to \infty} \frac{\sqrt{n}\mathbb{E}d_\infty\left(\bar{H}_n, C\right)}{\sqrt{2 \ln \ln n}} \leq \frac{1}{2}$ as well. Thus, the result follows. $\square$

**Theorem 3.2.** *For any copula $C$,*

$$\|C_{m,n} - C\|_p = O\left(\|B_m(C) - C\|_p + \left(\frac{m \ln \ln n}{n}\right)^{1/4p}\right) \text{ a.s.}$$

*and*

$$\mathbb{E}\|C_{m,n} - C\|_p = O\left(\|B_m(C) - C\|_p + \left(\frac{m \ln \ln n}{n}\right)^{1/4p}\right)$$

*as $m, n \to \infty$. In particular, $\|C_{m,n} - C\|_p \to 0$ whenever $m, n \to \infty$ with $m = o\left(\frac{n}{\ln \ln n}\right)$.*

***Proof.*** Since $0 \le \partial_u C_{m,n} \le 2$ and $0 \le \partial_u B_m(C) \le 1$, $|\partial_u C_{m,n} - \partial_u B_m(C)| \le 2$ and

$$\|C_{m,n} - B_m(C)\|_p^p \le \int_{\mathbb{I}} \int_{\mathbb{I}} 2^{p-1} |\partial_u C_{m,n}(u,v) - \partial_u B_m(C)(u,v)| \, du dv$$
$$= 2^{p-1} \|C_{m,n} - B_m(C)\|_1.$$

This implies

$$\|C_{m,n} - C\|_p = \|B_m(C) - C\|_p + \|C_{m,n} - B_m(C)\|_p$$
$$\le \|B_m(C) - C\|_p + 2^{(p-1)/p} \|C_{m,n} - B_m(C)\|_1^{1/p}.$$

Thus, it is sufficient to show that $\|C_{m,n} - B_m(C)\|_1 = O\left( \left(\frac{m \ln \ln n}{n}\right)^{1/4} \right)$ a.s. and $\mathbb{E} \|C_{m,n} - B_m(C)\|_1 = O\left( \left(\frac{m \ln \ln n}{n}\right)^{1/4} \right)$. Now,

$$|\partial_u C_{m,n} - \partial_u B_m(C)| \le \sum_{i=0}^m \sum_{j=0}^m \left| C_n\left(\frac{i}{m}, \frac{j}{m}\right) - C\left(\frac{i}{m}, \frac{j}{m}\right) \right| |\partial_n P_{i,m}(u)| P_{j,m}(v)$$
$$\le \sum_{i=0}^m \sum_{j=0}^m d_\infty(C_n, C) \left| \frac{k - mu}{u(1-u)} \right| P_{i,m}(u) P_{j,m}(v)$$
$$= \frac{d_\infty(C_n, C)}{u(1-u)} \sum_{i=0}^m |k - mu| P_{i,m}(u)$$
$$\le \frac{d_\infty(C_n, C)}{u(1-u)} \left( \sum_{i=0}^m |k - mu|^2 P_{i,m}(u) \right)^{1/2}$$
$$= \frac{d_\infty(C_n, C)}{u(1-u)} (mu(1-u))^{1/2}$$
$$= \frac{\sqrt{m} d_\infty(C_n, C)}{\sqrt{u(1-u)}}.$$

Thus,

$$\|C_{m,n} - B_m(C)\|_1 \le \sqrt{m} d_\infty(C_n, C) \int_{\mathbb{I}} \frac{1}{\sqrt{u(1-u)}} du$$
$$= \pi \sqrt{m} d_\infty(C_n, C)$$
$$= O\left( \sqrt{\frac{m \ln \ln n}{n}} \right)$$

a.s. as desire. Similarly,

$$\mathbb{E} \|C_{m,n} - B_m(C)\|_1 \le \pi \sqrt{m} \mathbb{E} d_\infty(C_n, C) = O\left( \sqrt{\frac{m \ln \ln n}{n}} \right).$$

$\square$

As an immediate application, we have the following result.

**Corollary 3.3.** *For any copula $C$, we have $\zeta_1\left(C_{m,n}\right) \to \zeta_1(C)$, $\omega\left(C_{m,n}\right) \to \omega(C)$, and $r\left(C_{m,n}\right) \to r(C)$ a.s. whenever $m, n \to \infty$ with $m = o\left(\frac{n}{\ln \ln n}\right)$. Moreover,*

$$\mathbb{E}\left(\zeta_1\left(C_{m,n}\right) - \zeta_1(C)\right)^2 = O\left(\left\|B_m(C) - C\right\|_1^2 + \left(\frac{m \ln \ln n}{n}\right)^{1/2}\right),$$

$$\mathbb{E}\left(r\left(C_{m,n}\right) - r(C)\right)^2 = O\left(\left\|B_m(C) - C\right\|_1^2 + \left(\frac{m \ln \ln n}{n}\right)^{1/2}\right), \quad and$$

$$\mathbb{E}\left(\omega\left(C_{m,n}\right) - \omega(C)\right)^2 = O\left(\left\|B_m(C) - C\right\|_1 + \left\|B_m(C^\perp) - C^\perp\right\|_1 + \left(\frac{m \ln \ln n}{n}\right)^{1/4}\right)$$

*where $C^\perp(u, v) = C(v, u)$ for all $u, v \in \mathbb{I}$.*

**Proof.** These statements follow from the fact that

$$\left|\zeta_1\left(C_{m,n}\right) - \zeta_1(C)\right| \le 3 \left\|C_{m,n} - C\right\|_1,$$

$$\left|r\left(C_{m,n}\right) - r(C)\right| \le 2 \left|\sqrt{r\left(C_{m,n}\right)} - \sqrt{r(C)}\right|$$

$$\le 12 \left\|C_{m,n} - C\right\|_2^2$$

$$\le 24 \left\|C_{m,n} - C\right\|_1,$$

and $\omega\left(C_{m,n}\right) = \sqrt{\frac{1}{2}\left(r\left(C_{m,n}\right) + r\left(C_{m,n}^\perp\right)\right)}$. $\qquad\qquad\square$

From the above result, in order to obtain the rate of convergence, we necessary have to obtain the rate that $B_m(C)$ converges to $C$ under Sobolev $L^1$-distance. We will provide an example of this computation in case the true copula $C$ has bounded second order derivatives. The proof is actually an adaptation of [16, Theorem 3].

**Theorem 3.4.** *Assume that the Hessian of a copula $C$ is bounded. Then,*

$$\left\|B_m(C) - C\right\|_1 = O\left(\frac{1}{\sqrt{m}}\right).$$

**Proof.** Consider the Taylor expansion,

$$C\left(\frac{i}{m}, \frac{j}{m}\right) = C(u, v) + \left(\frac{i}{m} - u\right) \partial_u C(u, v) + \left(\frac{j}{m} - v\right) \partial_v C(u, v)$$

$$+ \eta\left(\frac{i}{m}, u, \frac{j}{m}, v\right)\left(\left(\frac{i}{m} - u\right)^2 + \left(\frac{j}{m} - v\right)^2\right)$$

around the point $(u, v)$. Since the Hessian of the copula $C$ is bounded, there is a constant $M > 0$ such that $\left|\eta\left(\frac{i}{m}, u, \frac{j}{m}, v\right)\right| \le M$ for all $u, v \in \mathbb{I}$ and $i, j = 0, \ldots, m$. Now,

$$\partial_u B_m(C)(u, v) = \sum_{i=0}^{m} \sum_{j=0}^{m} C\left(\frac{i}{m}, \frac{j}{m}\right) \partial_u P_{i,m}(u) P_{j,m}(v)$$

$$= \sum_{i=0}^{m} \sum_{j=0}^{m} C(u, v) \partial_u P_{i,m}(u) P_{j,m}(v)$$

$$+ \sum_{i=0}^{m} \sum_{j=0}^{m} \left(\frac{i}{m} - u\right) \partial_u C(u, v) \partial_u P_{i,m}(u) P_{j,m}(v)$$

$$+ \sum_{i=0}^{m} \sum_{j=0}^{m} \left(\frac{j}{m} - v\right) \partial_v C(u, v) \partial_u P_{i,m}(u) P_{j,m}(v)$$

$$+ \sum_{i=0}^{m} \sum_{j=0}^{m} \xi\left(\frac{i}{m} - u, \frac{j}{m} - v\right) \partial_u P_{i,m}(u) P_{j,m}(v)$$

where $\xi\left(\frac{i}{m}-u, \frac{j}{m}-v\right) = \eta\left(\frac{i}{m}-u, \frac{j}{m}-v\right)\left(\left(\frac{i}{m}-u\right)^2+\left(\frac{j}{m}-v\right)^2\right)$. Notice that

$$\sum_{i=0}^{m}\sum_{j=0}^{m} C(u,v)\partial_u P_{i,m}(u)P_{j,m}(v) = \frac{C(u,v)}{u(1-u)}\sum_{i=0}^{m}(i-mu)P_{i,m}(u) = 0,$$

$$\sum_{i=0}^{m}\sum_{j=0}^{m}\left(\frac{i}{m}-u\right)\partial_u C(u,v)\partial_u P_{i,m}(u)P_{j,m}(v) = \frac{\partial_u C(u,v)}{mu(1-u)}\sum_{i=0}^{m}(i-mu)^2 P_{i,m}(u)$$
$$= \partial_u C(u,v),$$

and

$$\sum_{i=0}^{m}\sum_{j=0}^{m}\left(\frac{j}{m}-v\right)\partial_v C(u,v)\partial_u P_{i,m}(u)P_{j,m}(v)$$
$$= \frac{\partial_v C(u,v)}{mu(1-u)}\sum_{i=0}^{m}(i-mu)(j-mv)P_{i,m}(u)P_{j,m}(v)$$

$$\sum_{i=0}^{m}\sum_{j=0}^{m}\left(\frac{j}{m}-v\right)\partial_v C(u,v)\partial_u P_{i,m}(u)P_{j,m}(v)$$

$$= \frac{\partial_v C(u,v)}{mu(1-u)}\sum_{i=0}^{m}(i-mu)(j-mv)P_{i,m}(u)P_{j,m}(v)$$
$$= 0.$$

When $m \geq 2$, we have $3m-6 \geq 0$ and hence

$$|\partial_u B_m(C)(u,v) - \partial_u C(u,v)|$$
$$\leq \frac{mM}{u(1-u)}\sum_{i=0}^{m}\sum_{j=0}^{m}\left(\left(\frac{i}{m}-u\right)^2+\left(\frac{j}{m}-v\right)^2\right)\left|\frac{i}{m}-u\right|P_{i,m}(u)P_{j,m}(v)$$
$$\leq \frac{mM}{u(1-u)}\left(\sum_{i=0}^{m}\left(\frac{i}{m}-u\right)^4 P_{i,m}(u)\right)^{3/4}$$
$$\quad + \frac{mM}{u(1-u)}\sqrt{\sum_{i=0}^{m}\left(\frac{i}{m}-u\right)^2 P_{i,m}(u)\sum_{j=0}^{m}\left(\frac{j}{m}-v\right)^2 P_{j,m}(v)}$$
$$= \frac{mM}{u(1-u)}\left(\frac{mu(1-u)\left((3m-6)u(1-u)+1\right)}{m^4}\right)^{3/4}$$
$$\quad + \frac{mM}{u(1-u)}\left(\frac{v(1-v)}{m}\right)\sqrt{\frac{u(1-u)}{m}}$$
$$= \frac{M}{u^{1/4}(1-u)^{1/4}}\frac{\left((3m-6)u(1-u)+1\right)^{3/4}}{m^{5/4}} + \frac{Mv(1-v)}{\sqrt{mu(1-u)}}$$
$$\leq \frac{M}{u^{1/4}(1-u)^{1/4}}\frac{(3m/4)^{3/4}}{m^{5/4}} + \frac{M/4}{\sqrt{mu(1-u)}}$$
$$\leq \frac{3M/4}{\sqrt{mu(1-u)}} + \frac{M/4}{\sqrt{mu(1-u)}}$$
$$\leq \frac{M}{\sqrt{mu(1-u)}}$$

where the second inequality follows from Jensen's inequality and the third inequality follows from $x(1-x) \leq 1/4$ for all $x \in \mathbb{I}$. Thus,

$$\|B_m(C) - C\|_1 = \int_{\mathbb{I}} \frac{M}{\sqrt{mu(1-u)}} du = \frac{1}{\sqrt{m}} \pi M.$$

$\square$

**Remark 3.5.** The above result can also be relaxed to the case that the Hessian of $C$ is bounded only on $\Omega \subseteq \mathbb{I}^2$. In this case, $\|B_m(C) - C\|_1 = \int 1_{\mathbb{I}^2 \setminus \Omega} du dv + O\left(\frac{1}{\sqrt{m}}\right)$ which implies the same result as long as $\int 1_{\mathbb{I}^2 \setminus \Omega} du dv = O\left(\frac{1}{\sqrt{m}}\right)$. For example, in case of a Clayton copula $C_\theta(u,v) = \left(u^{-\theta} + v^{-\theta} - 1\right)^{-1/\theta}$ where $\theta > 0$, its Hessian is bounded on $\Omega_\delta = [\delta, 1]^2$ for all $\delta > 0$ but not for $\delta = 0$. Thus, $\|B_m(C_\theta) - C_\theta\|_1 = O\left(\frac{1}{\sqrt{m}}\right)$ also.

**Corollary 3.6.** *Assume that the Hessian of a copula $C$ is bounded, then we have*

$$\mathbb{E}\left(\zeta_1\left(C_{m,n}\right) - \zeta_1(C)\right)^2 = O\left(\frac{1}{m} + \left(\frac{m \ln \ln n}{n}\right)^{1/2}\right),$$

$$\mathbb{E}\left(r\left(C_{m,n}\right) - r(C)\right)^2 = O\left(\frac{1}{m} + \left(\frac{m \ln \ln n}{n}\right)^{1/2}\right), \quad and$$

$$\mathbb{E}\left(\omega\left(C_{m,n}\right) - \omega(C)\right)^2 = O\left(\frac{1}{\sqrt{m}} + \left(\frac{m \ln \ln n}{n}\right)^{1/4}\right).$$

Note that the above corollary can be used to find an optimal choice of $m$ for which the asymptotic mean square error of the estimator converges to zero the fastest. For example, consider the function $\phi(x) = \frac{1}{x} + \left(\frac{x \ln \ln n}{n}\right)^{1/2}$ where $0 \leq x \leq n$. This function has a unique minimum at $x = \left(\frac{4n}{\ln \ln n}\right)^{1/3}$. Thus, $m \approx 1.58 \left(\frac{n}{\ln \ln n}\right)^{1/3}$ will provide optimal rates of convergence for $\zeta_1\left(C_{m,n}\right)$ and $r\left(C_{m,n}\right)$ which are $O\left(\left(\frac{\ln \ln n}{n}\right)^{1/3}\right)$. Similarly, the function $\psi(x) = \frac{1}{\sqrt{x}} + \left(\frac{x \ln \ln n}{n}\right)^{1/4}$ has a unique minimum at $x = \left(\frac{16n}{\ln \ln n}\right)^{1/3}$. Thus, $m \approx 2.52 \left(\frac{n}{\ln \ln n}\right)^{1/3}$ will provide an optimal rate of convergence for $\omega\left(C_{m,n}\right)$ which is $O\left(\left(\frac{\ln \ln n}{n}\right)^{1/6}\right)$.

In general, we will also have

$$\|C_{m,n} - C\|_p = O\left(\left(\frac{\ln \ln n}{n}\right)^{1/6p}\right) \quad \text{a.s.}$$

and

$$\mathbb{E} \|C_{m,n} - C\|_p = O\left(\left(\frac{\ln \ln n}{n}\right)^{1/6p}\right)$$

where $m \approx \left(\frac{n}{\ln \ln n}\right)^{1/3}$ is an optimal choice for all $p \geq 1$. Note that this rate of convergence is slower than that in 2.2 and 3.1 which is normally expected.

Also, the kernel estimator for $r(C)$ has been provided in case the true copula $C$ is thrice continuously differentiable in the first variable and twice continuously differentiable in the second variable [3]. This assumption is stronger than requiring the Hessian of $C$ to be bounded. Nevertheless, the result in [3] is also stronger than ours. In [3], Dette et al. were able to prove the central limit theorem for their estimator while our result is only that of the law of large number.

## 4. Numerical study

### 4.1. Simulation

In this part, we will study asymptotic behavior of the estimators $\zeta_1(C_{m,n})$, $\omega(C_{m,n})$, and $r(C_{m,n})$ via simulations. We already proved that these estimators converge to their respective measures $\zeta_1(C)$, $\omega(C)$, and $r(C)$. Rate of convergences are also given up to some constants which might be useless if these constants are too large. We will confirm that this is not the case by showing that the asymptotic errors is already acceptable for a reasonable sample size $n = 50, 100, 200$, and $400$.

The value $m = m(n)$ is chosen to optimize the rate of convergence as stated at the end of the previous section. Specifically, $m \approx 1.58\left(\frac{n}{\ln\ln n}\right)^{1/3}$ for $\zeta_1(C_{m,n})$ and $r(C_{m,n})$ while $m \approx 2.52\left(\frac{n}{\ln\ln n}\right)^{1/3}$ for $\omega(C_{m,n})$. Thus, the pair $(n,m)$ of the sample size $n$ and the value $m = m(n)$ will be $(50,5)$, $(100,6)$, $(200,8)$, and $(400,10)$ for $\zeta_1(C_{m,n})$ and $r(C_{m,n})$ while $(n,m)$ will be $(50,8)$, $(100,10)$, $(200,13)$ and $(400,15)$ for $\omega(C_{m,n})$. Asymptotic mean square error (AMSE), asymptotic bias, and asymptotic variance based on $10,000$ simulation runs will be given. The families of Clayton copulas will be chosen as true copulas for the simulation. The result of the simulation is as follows.

**Table 2.** The AMSE of $\zeta_1(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $8.05 \times 10^{-3}$ | $8.69 \times 10^{-3}$ | $2.29 \times 10^{-2}$ | $4.94 \times 10^{-2}$ |
| 100 | $5.37 \times 10^{-3}$ | $7.14 \times 10^{-3}$ | $1.80 \times 10^{-2}$ | $3.88 \times 10^{-2}$ |
| 200 | $3.57 \times 10^{-3}$ | $4.13 \times 10^{-3}$ | $9.90 \times 10^{-3}$ | $2.22 \times 10^{-2}$ |
| 400 | $2.25 \times 10^{-3}$ | $2.69 \times 10^{-3}$ | $6.62 \times 10^{-3}$ | $1.52 \times 10^{-2}$ |

**Table 3.** The asymptotic bias of $\zeta_1(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $8.21 \times 10^{-2}$ | $-7.33 \times 10^{-2}$ | $-1.39 \times 10^{-1}$ | $-2.17 \times 10^{-1}$ |
| 100 | $6.82 \times 10^{-2}$ | $-7.07 \times 10^{-2}$ | $-1.26 \times 10^{-1}$ | $-1.94 \times 10^{-1}$ |
| 200 | $5.61 \times 10^{-2}$ | $-5.20 \times 10^{-2}$ | $-9.29 \times 10^{-2}$ | $-1.46 \times 10^{-1}$ |
| 400 | $4.51 \times 10^{-2}$ | $-4.29 \times 10^{-2}$ | $-7.67 \times 10^{-2}$ | $-1.21 \times 10^{-1}$ |

**Table 4.** The asymptotic variance of $\zeta_1(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $1.30 \times 10^{-3}$ | $3.31 \times 10^{-3}$ | $3.56 \times 10^{-3}$ | $2.36 \times 10^{-3}$ |
| 100 | $7.29 \times 10^{-4}$ | $2.14 \times 10^{-3}$ | $2.06 \times 10^{-3}$ | $1.30 \times 10^{-3}$ |
| 200 | $4.17 \times 10^{-4}$ | $1.43 \times 10^{-3}$ | $1.27 \times 10^{-3}$ | $8.04 \times 10^{-4}$ |
| 400 | $2.17 \times 10^{-4}$ | $8.46 \times 10^{-4}$ | $7.34 \times 10^{-4}$ | $4.62 \times 10^{-4}$ |

First, consider the case of $\zeta_1(C_{m,n})$. The asymptotic mean square of the simulation seem acceptable with only a small contribution form the asymptotic variance. The bias, however, is a bit high when $\theta \geq 1$. Overall, the estimator seem to perform worse when $\theta$ becomes bigger suggesting that larger sample size is needed for a more accurate estimation.

**Table 5.** The AMSE of $\omega(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $4.62 \times 10^{-2}$ | $2.73 \times 10^{-3}$ | $1.75 \times 10^{-2}$ | $3.98 \times 10^{-2}$ |
| 100 | $7.28 \times 10^{-2}$ | $2.83 \times 10^{-3}$ | $2.89 \times 10^{-3}$ | $1.25 \times 10^{-2}$ |
| 200 | $5.90 \times 10^{-2}$ | $1.09 \times 10^{-3}$ | $3.24 \times 10^{-3}$ | $1.11 \times 10^{-2}$ |
| 400 | $6.27 \times 10^{-2}$ | $1.67 \times 10^{-3}$ | $1.35 \times 10^{-3}$ | $6.55 \times 10^{-3}$ |

**Table 6.** The asymptotic bias of $\omega(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $2.11 \times 10^{-1}$ | $-3.14 \times 10^{-2}$ | $-1.21 \times 10^{-1}$ | $-1.94 \times 10^{-1}$ |
| 100 | $2.69 \times 10^{-1}$ | $4.44 \times 10^{-2}$ | $-3.82 \times 10^{-2}$ | $-1.06 \times 10^{-1}$ |
| 200 | $2.43 \times 10^{-1}$ | $2.20 \times 10^{-2}$ | $-4.81 \times 10^{-2}$ | $-1.02 \times 10^{-1}$ |
| 400 | $2.50 \times 10^{-1}$ | $3.65 \times 10^{-2}$ | $-2.92 \times 10^{-2}$ | $-7.86 \times 10^{-2}$ |

**Table 7.** The asymptotic variance of $\omega(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $1.51 \times 10^{-3}$ | $1.75 \times 10^{-3}$ | $2.82 \times 10^{-3}$ | $2.36 \times 10^{-3}$ |
| 100 | $2.76 \times 10^{-4}$ | $8.55 \times 10^{-4}$ | $1.42 \times 10^{-3}$ | $1.18 \times 10^{-3}$ |
| 200 | $2.01 \times 10^{-4}$ | $6.01 \times 10^{-4}$ | $9.27 \times 10^{-4}$ | $7.20 \times 10^{-4}$ |
| 400 | $7.77 \times 10^{-5}$ | $3.34 \times 10^{-4}$ | $5.00 \times 10^{-4}$ | $3.68 \times 10^{-4}$ |

Next, consider the case of $\omega(C_{m,n})$. Again, most of the contribution to the AMSE is from the asymptotic bias. When comparing to $\zeta_1(C_{m,n})$, the estimator seem to perform worse when $\theta = 0$. The situation is a bit different with $r(C_{m,n})$. In this case, the AMSE, the asymptotic bias, and variance all seem reasonable. Comparing to the kernel estimator by Dette et al.[3], both AMSE are roughly of the same level, the asymptotic bias of the Bernstein estimator are higher but the asymptotic variance are significantly lower.

**Table 8.** The AMSE of $r(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $5.16 \times 10^{-3}$ | $4.08 \times 10^{-4}$ | $4.45 \times 10^{-3}$ | $3.39 \times 10^{-2}$ |
| 100 | $2.99 \times 10^{-3}$ | $1.54 \times 10^{-4}$ | $5.87 \times 10^{-3}$ | $3.44 \times 10^{-2}$ |
| 200 | $4.37 \times 10^{-3}$ | $3.97 \times 10^{-4}$ | $2.15 \times 10^{-3}$ | $1.82 \times 10^{-2}$ |
| 400 | $4.17 \times 10^{-3}$ | $4.28 \times 10^{-4}$ | $1.35 \times 10^{-3}$ | $1.29 \times 10^{-2}$ |

**Table 9.** The asymptotic bias of $r(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $7.11 \times 10^{-2}$ | $1.32 \times 10^{-2}$ | $-6.24 \times 10^{-2}$ | $-1.82 \times 10^{-1}$ |
| 100 | $5.41 \times 10^{-2}$ | $-3.04 \times 10^{-3}$ | $-7.42 \times 10^{-2}$ | $-1.84 \times 10^{-1}$ |
| 200 | $6.59 \times 10^{-2}$ | $1.64 \times 10^{-2}$ | $-4.28 \times 10^{-2}$ | $-1.33 \times 10^{-1}$ |
| 400 | $6.44 \times 10^{-2}$ | $1.85 \times 10^{-2}$ | $-3.37 \times 10^{-2}$ | $-1.12 \times 10^{-1}$ |

**Table 10.** The asymptotic variance of $r(C_{m,n})$ for Clayton copulas $C_\theta$.

| $n\backslash\theta$ | 0.0 | 0.5 | 1.0 | 2.0 |
|---|---|---|---|---|
| 50 | $1.12 \times 10^{-4}$ | $1.75 \times 10^{-4}$ | $5.59 \times 10^{-4}$ | $7.92 \times 10^{-4}$ |
| 100 | $6.62 \times 10^{-5}$ | $8.55 \times 10^{-4}$ | $3.70 \times 10^{-4}$ | $5.27 \times 10^{-4}$ |
| 200 | $3.04 \times 10^{-5}$ | $6.01 \times 10^{-4}$ | $3.16 \times 10^{-4}$ | $4.27 \times 10^{-4}$ |
| 400 | $1.68 \times 10^{-5}$ | $3.34 \times 10^{-4}$ | $2.11 \times 10^{-4}$ | $2.69 \times 10^{-4}$ |

## 4.2. Data Sample

In this part, we provide an example study of a relationship between dust density, temperature, and humidity. The data is taken from the Climate Change Data Center of Chiang Mai University[‡]. The hourly average data from 6 censors within the Mueang District between September 14th, 2018 and September 25th, 2018 with the total of 966 sets of data are used to compute the value of $\zeta_1(Y|X)$, $r(Y|X)$, and $\omega(X,Y)$ where $X$ is either the temperature or the humidity level and $Y$ is either the density $(ug/m^3)$ of dust with diameter at most 10 micron (PM10) or the density of dust with diameter at most 2.5 micron (PM2.5). The Pearson correlation $\text{Corr}(X,Y)$ and the Spearman's rank correlation $\rho(X,Y)$ between $X$ and $Y$ are also computed for comparison.

**Table 11.** Estimate values for assoication level among climate data within the Mueang District, Chiang Mai, Thailand

| $X$ | Temperature | | Humidity level | |
|---|---|---|---|---|
| $Y$ | PM2.5 density | PM10 density | PM2.5 density | PM10 density |
| $\zeta_1(Y|X)$ | 0.16386 | 0.16678 | 0.44264 | 0.42667 |
| $r(Y|X)$ | 0.08887 | 0.09397 | 0.32856 | 0.31442 |
| $\omega(X,Y)$ | 0.29978 | 0.30255 | 0.67051 | 0.65936 |
| $\text{Corr}(X,Y)$ | -0.01068 | -0.04218 | 0.02354 | 0.08263 |
| $\rho(X,Y)$ | 0.09650 | 0.06329 | -0.06942 | -0.00724 |

From the estimated value, the Pearson correlation between $X$ and $Y$ are small suggesting that they do not have a linear relationship. The Spearman's rank correlation is also quite small which again suggesting that $X$ and $Y$ might not have monotone relationship. The higher values of $\zeta_1(Y|X)$, $r(Y|X)$, and $\omega(X,Y)$ suggest that the dust density actually depends on the temperature and the humidity level to a certain degree. Their relationship, however, might be nonlinear and might not be monotone. The fact that $\zeta_1(Y|X)$, $r(Y|X)$, and $\omega(X,Y)$ are higher when $X$ is the humidity level also suggest that the dust density depends more on humidity level than the temperature. The fact that these values are not close to one is only natural since there are more than one variables (temperature and humidity level) that are the source of the dust density. Perhaps, an estimator of a multivariate version of these measures have to be constructed and used in the studied of these data in the future.

---

[‡]website: www.cmuccdc.org

# References

[1] J. Behboodian, A. Dolati, and M. Úbeda-Flores. A multivariate version of gini's rank association coefficient. *Statist. Papers*, 48(2):295–304, 2007.

[2] N. Blomqvist. On a measure of dependence between two random variables. *Ann. Math. Statist.*, 21(4):593–600, 1950.

[3] H. Dette, K. F. Siburg, and P. A. Stoimenov. A copula-based non-parametric measure of regression dependence. *Scand. J. Stat.*, 40(1):21–41, 2013.

[4] S. Gaißer, M. Ruppert, and F. Schmid. A multivariate version of hoeffding's phi-square. *J. Multivariate Anal.*, 101(10):2571–2586, 2010.

[5] W. Hoeffding. *The collected works of Wassily Hoeffding.* Springer-Verlag, 1994.

[6] P. Janssen, J. Swanepoel, and N. Veraverbeke. Large sample behavior of the bernstein copula estimator. *J. Statist. Plann. Inference*, 142(5):1189 – 1197, 2012.

[7] H. Joe. Multivariate concordance. *J. Multivariate Anal.*, 35(1):12–30, 1990.

[8] M. G. Kendall. A new measure of rank correlation. *Biometrika*, 30(1/2):81–93, 1938.

[9] H. O. Lancaster. Measures and indeces of dependence. In M. Kotz and N. L. Johnson, editors, *Encyclopedia of Statistical Sciences*, volume 2, pages 334–339. Wiley, New York, 1982.

[10] R. B. Nelsen. Nonparametric measures of multivariate association. *Lecture Notes-Monograph Series*, 28:223–232, 1996.

[11] A. Rényi. On measures of dependence. *Acta Math. Hungar.*, 10(3):441–451, 1959.

[12] F. Schmid, R. Schmidt, T. Blumentritt, S. Gaißer, and M. Ruppert. Copula-based measures of multivariate association. In P. Jaworski, F. Durante, W. K. Härdle, and T. Rychlik, editors, *Copula Theory and Its Applications*, volume 198 of *Lecture Notes in Statistics – Proceedings*, pages 209–236. Springer Berlin Heidelberg, 2010.

[13] B. Schweizer and E. F. Wolff. On nonparametric measures of dependence for random variables. *Ann. Statist.*, 9(4):879–885, 1981.

[14] K. F. Siburg and P. A. Stoimenov. A measure of mutual complete dependence. *Metrika*, 71:239–251, 2010.

[15] A. Sklar. Fonctions de répartition á n dimensions et leurs marges. *Publ. Inst. Statist. Univ. Paris*, 8:229–231, 1959.

[16] M. Taylor. Bernstein Polynomials and n-Copulas. *ArXiv e-prints*, March 2009.

[17] M. D. Taylor. Multivariate measures of concordance. *Ann. Inst. Statist. Math.*, 59(4):789–806, 2006.

[18] W. Trutsching. On a strong metric on the space of copulas and its induced dependence measure. *J. Math. Anal. Appl.*, 384:690–705, 2011.

[19] M. Úbeda-Flores. Multivariate versions of blomqvist's beta and spearman's footrule. *Ann. Inst. Statist. Math.*, 57(4):781–788, 2005.

[20] E. F. Wolff. N-dimensional measures of dependence. *Stochastica*, 4(3):175–188, 1980.