

Multivariate Copula Expressed by Lower Dimensional Copulas

Edith Kovács^a, Tamás Szántai^b

^a*Department of Mathematics, ÁVF College of Management of Budapest, Villányi út
11-13, H-1114 Budapest, Hungary*

Corresponding author,

E-mail address: kovacs.edith@avf.hu

^b*Institute of Mathematics, Budapest University of Technology and Economics,
Műegyetem rkp. 3, H-1111 Budapest, Hungary*

E-mail address: szantai@math.bme.hu

Abstract

Modeling of high order multivariate probability distribution is a difficult problem which occurs in many fields. Copula approach is a good choice for this purpose, but the curse of dimensionality still remains a problem. In this paper we give a theorem which expresses a multivariate copula by using only some lower dimensional ones based on the conditional independences between the variables. In general the construction of a multivariate copula using this theorem is quite difficult, due the consistency properties which have to be fulfilled. For this purpose we introduce the sample derived copula, and prove that the dependence between the random variables involved depends just on this copula and on the partition. By using the sample derived copula the theorem can be successfully applied, in order to construct a multivariate discrete copula by using some of its marginals.

Keywords:

Multivariate copula, Junction tree, Conditional independence, Sample derived copula.

1. Introduction

First we motivate why should we model the multivariate distribution by copulas from an information theoretical point of view. The information content of a multivariate probability distribution depends only on its copula

density. In [12] and [1] one can see this result for the two-dimensional case and the same is true for more dimensions, too.

In this paper we prove a theorem which links the multivariate probability distribution assigned to a junction tree to the multivariate copula. It is known that the probability distribution assigned to a junction tree uses the conditional independence structure underlying the random variables so the copula introduced here will have this property, too.

In this introductory part we describe the main concepts and introduce the notations which we will use in the paper. In the second section we prove a theorem which links a multivariate copula to the junction tree probability distribution. In the third section we will introduce the concept of Sample Derivated Copula (SDC) which makes possible the exploitation of the conditional independences between the random variables. We prove that the information content of the probability distribution given by a partition set depends only on the SDC. In the fourth section we apply the junction tree approach to the SDC.

We finish the paper with conclusions and possible applications.

Let $V = \{1, \dots, n\}$ be a set of vertices. A hypergraph is a set V of vertices together with a set Γ of subsets of V . A hypergraph is acyclic if no elements in Γ are subsets of other elements, and if the elements of Γ can be ordered (K_1, \dots, K_m) to have the *running intersection property*: for all $j \geq 2$, exists $i < j : K_i \supseteq K_j \cap (K_1 \cup \dots \cup K_{j-1})$ [8].

It is convenient to introduce the so called separator sets $S_j = K_j \cap (K_1 \cup \dots \cup K_{j-1})$, where $S_1 = \phi$.

We note here that if $R_j = K_j \setminus S_j$ then S_j separates (in graph terms) the vertices in R_j from the vertices in $(K_1 \cup \dots \cup K_{j-1}) \setminus S_j$.

We mention here that a hypergraph (V, Γ) is acyclic if and only if Γ can be considered to be the set of cliques of a chordal (triangulated) graph [9],[16].

In the following we consider acyclic hypergraphs with the property that the union of all sets in Γ is V . We denote the separator set by \mathcal{S} and refer to the acyclic hypergraph as (V, Γ, \mathcal{S}) .

Let $V = \{1, 2, \dots, n\}$ be the set of indices of the continuous random variables $X = \{X_1, \dots, X_n\}$. We suppose that the probability density functions of X_1, \dots, X_n exist and denote them by f_{X_1}, \dots, f_{X_n} .

We need the following notations:

- $F_{X_i}(x_i) = P(X_i < x_i; X_j = \infty \text{ for all } j \neq i)$ stands for the univariate marginal cumulative distribution function corresponding to the variable

X_i ,

- The joint probability density function and the joint cumulative distribution function of $(X_1, \dots, X_n)^T$ is denoted by $f_{\mathbf{X}}(\mathbf{x})$ and $F_{\mathbf{X}}(\mathbf{x})$, respectively,
- $D = \{i_1, \dots, i_d\} \subset V$, $\mathbf{X}_D = (X_{i_1}, \dots, X_{i_d})^T$, $\mathbf{x}_D = (x_{i_1}, \dots, x_{i_d})^T$,
- The d -th order marginal probability density function and the d -th order marginal cumulative distribution function of \mathbf{X}_D is denoted by $f_{\mathbf{X}_D}(\mathbf{x}_D)$ and $F_{\mathbf{X}_D}(\mathbf{x}_D)$, respectively.

Having these notations we give the concept of the junction tree. It is known that the junction tree encodes the conditional independences between the variables. Let us remark here that from now on the indices of the random variables are assigned to the nodes of a graph. In the graph a set of nodes B separates a set of nodes A from another set of nodes C , where A, B, C are disjoint subsets of V , if and only if X_A and X_C are conditionally independent with respect to X_B (see the definition of the Markov random field).

Definition 1. A junction tree over X is a cluster tree, which is assigned to an acyclic hypergraph (V, Γ, \mathcal{S}) as follows:

- 1) Each cluster of the cluster tree consists of a subset X_K of X , where $K \in \Gamma$. To each cluster is assigned the joint marginal density function $f_{\mathbf{X}_K}(\mathbf{x}_K)$;
- 2) Each edge connecting to clusters is called separator and consists of a subset X_S of X , where S is a separator set. To each separator there is assigned the marginal probability density function $f_{\mathbf{X}_S}(\mathbf{x}_S)$;
- 3) The union of all clusters is X .

Definition 2. A junction tree probability distribution is a probability distribution assigned to the junction tree in the following way:

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{\prod_{K \in \Gamma} f_{\mathbf{X}_K}(\mathbf{x}_K)}{\prod_{S \in \mathcal{S}} \left(f_{\mathbf{X}_S}(\mathbf{x}_S) \right)^{v_S - 1}},$$

where v_S is the number of those clusters which contain all the variables of \mathbf{X}_S .

It is useful to note here that since in the hypergraph (V, Γ, \mathcal{S}) S_j separates (in graph terms) the vertices in $R_j = K_j - S_j$ from the vertices in $(K_1 \cup \dots \cup K_{j-1}) - S_j$ the random variables with indices in $R_j = K_j - S_j$ and the variables with indices in $(K_1 \cup \dots \cup K_{j-1}) - S_j$ are conditionally independent with respect to the variables with indices in S_j .

Remark 1. Since the junction tree is assigned to an acyclic hypergraph, the running intersection property stands for the junction tree, too. It can be reformulated as follows. If two clusters contain a random variable, then all clusters on the path between these clusters contain this random variable.

First we call back the concept of copula and formulate the Sklar's theorem (see [3] and [13]).

Definition 3. A function $C : [0; 1]^d \rightarrow [0; 1]$ is called a d -dimensional copula if it satisfies the following conditions:

- 1) $C(u_1, \dots, u_d)$ is increasing in each component u_i ,
- 2) $C(u_1, \dots, u_{i-1}, 0, u_{i+1}, \dots, u_d) = 0$ for all $u_k \in [0; 1]$, $k \neq i$, $i = 1, \dots, n$,
- 3) $C(1, \dots, 1, u_i, 1, \dots, 1) = u_i$ for all $u_i \in [0; 1]$, $i = 1, \dots, d$,
- 4) C is d -increasing, i.e for all $(u_{1,1}, \dots, u_{1,d})$ and $(u_{2,1}, \dots, u_{2,d})$ in $[0; 1]^d$ with $u_{1,i} < u_{2,i}$ for all i , we have

$$\sum_{i_1=1}^2 \dots \sum_{i_d=1}^2 (-1)^{\sum_{j=1}^d i_j} C(u_{i_1,1}, \dots, u_{i_d,d}) \geq 0.$$

Due to Sklar's theorem if X_1, \dots, X_d are continuous random variables defined on a common probability space, with the univariate marginal cdf's $F_{X_i}(x_i)$ and the joint cdf $F_{X_1, \dots, X_d}(x_1, \dots, x_d)$ then there exists a unique copula function $C_{X_1, \dots, X_d}(u_1, \dots, u_d) : [0; 1]^d \rightarrow [0; 1]$ such that by the substitution $u_i = F_i(x_i)$, $i = 1, \dots, d$ we get

$$F_{X_1, \dots, X_d}(x_1, \dots, x_d) = C_{X_1, \dots, X_d}(F_1(x_1), \dots, F_d(x_d))$$

for all $(x_1, \dots, x_d)^T \in R^d$.

In the following we will use the vectorial notation $F_{\mathbf{X}_D}(\mathbf{x}_D) = C_{X_D}(\mathbf{u}_D)$, where $\mathbf{u}_D = (F_{X_{i_1}}(x_{i_1}), \dots, F_{X_{i_d}}(x_{i_d}))^T$.

We need the following assertion:

$$\begin{aligned}
& f_{X_{i_1}, \dots, X_{i_d}}(x_{i_1}, \dots, x_{i_d}) = \\
&= \frac{\partial^d F_{X_{i_1}, \dots, X_{i_d}}(x_{i_1}, \dots, x_{i_d})}{\partial x_{i_1} \cdots \partial x_{i_d}} \\
&= \frac{\partial^d C_{X_{i_1}, \dots, X_{i_d}}(F_{X_{i_1}}(x_{i_1}), \dots, F_{X_{i_d}}(x_{i_d}))}{\partial x_{i_1} \cdots \partial x_{i_d}} \\
&= \frac{\partial^d C_{X_{i_1}, \dots, X_{i_d}}(u_{i_1}, \dots, u_{i_d})}{\partial u_{i_1} \cdots \partial u_{i_d}} \Bigg|_{u_{i_k} = F_{X_{i_k}}(x_{i_k}), k=1, \dots, d} \cdot \prod_{k=1}^d \frac{\partial F_{X_{i_k}}(x_{i_k})}{\partial x_{i_k}} \\
&= c_{X_{i_1}, \dots, X_{i_d}}(F_{X_{i_1}}(x_{i_1}), \dots, F_{X_{i_d}}(x_{i_d})) \cdot \prod_{k=1}^d f_{X_{i_k}}(x_{i_k})
\end{aligned}$$

In vectorial notation this can be written as

$$f_{\mathbf{X}_D}(\mathbf{x}_D) = c_{\mathbf{X}_D}(\mathbf{u}_D) \cdot \prod_{i_k \in D} f_{X_{i_k}}(x_{i_k}) \quad (1)$$

and from (1) we get

$$c_{\mathbf{X}_D}(\mathbf{u}_D) = \frac{f_{\mathbf{X}_D}(\mathbf{x}_D)}{\prod_{i_k \in D} f_{X_{i_k}}(x_{i_k})} \quad (2)$$

2. The multivariate copula associated to a junction tree probability distribution.

Theorem 1. *The copula density function associated to a junction tree probability distribution*

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{\prod_{K \in \Gamma} f_{\mathbf{X}_K}(\mathbf{x}_K)}{\prod_{S \in \mathcal{S}} [f_{\mathbf{X}_S}(\mathbf{x}_S)]^{v_S - 1}},$$

is given by

$$c_{\mathbf{X}}(\mathbf{u}_V) = \frac{\prod_{K \in \Gamma} c_{\mathbf{X}_K}(\mathbf{u}_K)}{\prod_{S \in \mathcal{S}} [c_{\mathbf{X}_S}(\mathbf{u}_S)]^{v_S - 1}}. \quad (3)$$

PROOF.

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{\prod_{K \in \Gamma} f_{\mathbf{X}_K}(\mathbf{x}_K)}{\prod_{S \in \mathcal{S}} [f_{\mathbf{X}_S}(\mathbf{x}_S)]^{v_S-1}} = \frac{\prod_{K \in \Gamma} c_{\mathbf{X}_K}(\mathbf{u}_K) \cdot \prod_{i_k \in K} f_{X_{i_k}}(x_{i_k})}{\prod_{S \in \mathcal{S}} \left[c_{\mathbf{X}_S}(\mathbf{u}_S) \cdot \prod_{i_k \in S} f_{X_{i_k}}(x_{i_k}) \right]^{v_S-1}}. \quad (4)$$

The question that we have to answer is how many times appears in the nominator respectively in the denominator the probability density function $f_{X_i}(x_i)$ of each X_i random variable.

Since $\bigcup_{K \in \Gamma} \mathbf{X}_K = X$ for each random variable $X_i \in X$, $f_{X_i}(x_i)$ appears at least once in the nominator.

Now we prove that in the junction tree over X the number of clusters which contain a variable X_i is greater with 1 than the number of separators which contain the same variable. This is true for all $i = 1, \dots, n$. This means $\#\{K \in \Gamma | X_i \in X_K\} = \#\{S \in \mathcal{S} | X_i \in X_S\} + 1$.

For a variable X_i we denote $\#\{S \in \mathcal{S} | X_i \in X_S\}$ by t .

Case: $t = 0$.

The statement is a consequence of the definition of junction tree, that is the union of all clusters is X , so every variable have to appear at least in one cluster. X_i can not appear in two clusters, because in this case there should exist a separator which contain X_i too, and we supposed that there is not such a separator ($t = 0$)

Case: $t > 0$

If two clusters contain the variable X_i , then every cluster from the path between the two clusters contain X_i (running intersection property). From this results that the clusters containing X_i are the nodes of a connected graph, and this graph is a tree. If this tree contain t separator sets then it contains $t + 1$ clusters. All of these separators contain X_i , and each separator connects two clusters. So there will be $t + 1$ clusters that contain X_i .

Applying this result in formula (4) after simplification we obtain

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{\prod_{K \in \Gamma} c_{\mathbf{X}_K}(\mathbf{u}_K) \prod_{i=1}^n f_{X_i}(x_i)}{\prod_{S \in \mathcal{S}} [c_{\mathbf{X}_S}(\mathbf{u}_S)]^{v_S-1}}.$$

Dividing both sides by $\prod_{i=1}^n f_{X_i}(x_i)$ we obtain

$$\frac{f_{\mathbf{X}}(\mathbf{x})}{\prod_{i=1}^n f_{X_i}(x_i)} = \frac{\prod_{K \in \Gamma} c_{\mathbf{X}_K}(\mathbf{u}_K)}{\prod_{S \in \mathcal{S}} [c_{\mathbf{X}_S}(\mathbf{u}_S)]^{v_S-1}}. \quad (5)$$

Equations (2) and (5) prove the statement of the theorem.

We saw that if the conditional independence structure underlying the random variables makes possible the construction of a junction tree, then the multivariate copula density associated to the joint probability distribution can be expressed as a product and fraction of lower dimensional copula densities.

A logical question is the following. What conditions are necessary for (5) to be a copula density? It is easy to see that the product and fraction of copulas are positive. So

$$c(\mathbf{u}) = \frac{\prod_{K \in \Gamma} c_K(\mathbf{u}_K)}{\prod_{S \in \mathcal{S}} [c_S(\mathbf{u}_S)]^{v_S-1}}$$

will be a copula density if and only if

$$\int_{[0;1]^n} \frac{\prod_{K \in \Gamma} c_K(\mathbf{u}_K)}{\prod_{S \in \mathcal{S}} [c_S(\mathbf{u}_S)]^{v_S-1}} d\mathbf{u} = 1.$$

This happens if the following consistency conditions are fulfilled for all connected clique pairs K_i and K_j :

$$\int_{[0;1]^{\#\{K_i \setminus S_{ij}\}}} c_{K_i}(\mathbf{u}_{K_i}) d\mathbf{u}_{K_i \setminus S_{ij}} = \int_{[0;1]^{\#\{K_j \setminus S_{ij}\}}} c_{K_j}(\mathbf{u}_{K_j}) d\mathbf{u}_{K_j \setminus S_{ij}},$$

where $S_{ij} = K_i \cap K_j$. We emphasize here that all cliques are subsets of the set Γ of the acyclic hypergraph (V, Γ, \mathcal{S}) .

These conditions are fulfilled if $c_{S_{ij}}(\mathbf{u}_{S_{ij}})$ are marginal probability densities of $c_{K_i}(\mathbf{u}_{K_i})$, whenever S_{ij} connects a cluster K_i . This can be expressed by terms of copula function as follows.

For $\{k_1, \dots, k_m\} = K_i$ and $\{s_1, \dots, s_l\} = S_{ij}$, $\{s_1, \dots, s_l\} \subset \{k_1, \dots, k_m\}$ stands $C_m(u_{k_1}, \dots, u_{k_m}) = C_l(u_{s_1}, \dots, u_{s_l})$ for $u_{k_i} = 1$, when $k_i \notin S_{ij}$. Usually this condition is not fulfilled by copulas.

Finding multivariate copulas which fulfill the consistency conditions is not a trivial task.

A special type of conditional independence, when the graph underlying the random variables is starlike, is treated in [19]. Another type of special multivariate copula where the underlying conditional independence graph is a tree can be found in [6].

For discrete random variables, the conditional independences are exploited by the Markov random fields. In physics for two-valued random variables it is known the Ising model. In these cases the random variables take on a few values only. However many times the problem is hard. The great advantage of using the discrete approach is that the marginal probability distributions involved fulfill the consistency conditions (see [17] and [18]).

If we have an i.i.d. sample of size N from a continuous joint probability distribution then for each random variable we have N different values. For this case, the empirical copulas were introduced and first studied by P. Deheuvels in [5] who called them empirical dependence functions. Later in [10] and [11] there were introduced the so called discrete copulas. About the two-dimensional empirical copulas one can read in Nelsen's introductory book (see [13]). In the case when we are dealing with a sample drawn from a continuous joint probability distribution the size of these random variables would be too large, so we will apply a uniform partition and define the so called sample derivated copula.

3. The sample derivated copula.

Let X_1, \dots, X_n be continuous random variables in the same probability field. Let

$$\begin{array}{c} x_1^1, \dots, x_n^1 \\ x_1^2, \dots, x_n^2 \\ \vdots \\ x_1^N, \dots, x_n^N \end{array} \tag{6}$$

be an i.i.d. sample of size N taken from the joint probability distribution of the random vector $(X_1, \dots, X_n)^T$.

As any sample element occurs two times in the sample with probability zero, we can suppose that the sample elements are different.

We denote the set of the values of X_i in the sample by Λ_i . This set contains N values, for each random variable. The theoretical range of the continuous random variable X_i will be denoted by $\overline{\Lambda}_i$. For every i we denote by $\lambda_i^m = \min \overline{\Lambda}_i \in R$ and by $\lambda_i^M = \max \overline{\Lambda}_i \in R$. We suppose for simplicity that $\min \overline{\Lambda}_i \neq \min \Lambda_i$ and $\max \overline{\Lambda}_i \neq \max \Lambda_i$. For each random variable X_i we define a partition of Λ_i by $\mathcal{P}_i = \{x_0^{p_i} = \lambda_i^m, x_1^{p_i}, \dots, x_{m_i-1}^{p_i}, x_{m_i}^{p_i} = \lambda_i^M\}$ with the following properties:

- For each random variable X_i , each interval $(x_{j-1}^{p_i}; x_j^{p_i}]$, $j = 1, \dots, m_i$ contains a given $n_i = \frac{N}{m_i} \in N$ number of values from the set Λ_i .
- Each $x_j^{p_i} \in \Lambda_i$, $j = 1, \dots, m_i - 1$.

The partition with the above properties will be called uniform partition. We denote by \mathcal{P} the set of partitions $\{\mathcal{P}_1, \dots, \mathcal{P}_n\}$.

Let be \widetilde{X}_i the categorical random variable associated to the random variable X_i :

$$P(\widetilde{X}_i \in (x_{j-1}^{p_i}; x_j^{p_i}]) = \frac{1}{m_i}, j = 1, \dots, m_i.$$

We assign to each $x^i \in (x_{j-1}^{p_i}; x_j^{p_i}]$ the number $u_j^i = \frac{j}{m_i}$, $j = 0, \dots, m_i$.

Obviously $u_0^i = 0$ and $u_{m_i}^i = 1$. Let $\widetilde{\Lambda}_i = \{u_j^i | j = 0, \dots, m_i\}$. So we can define the following discrete uniform random variables:

$$\widetilde{U}_i = \begin{pmatrix} u_0^i & u_1^i & \dots & u_{m_i-1}^i & u_{m_i}^i \\ 0 & \frac{1}{m_i} & \dots & \frac{1}{m_i} & \frac{1}{m_i} \end{pmatrix}, i = 1, \dots, n.$$

Now we transform the sample (6) using the above assignment. We denote the transformed sample by \mathcal{T} .

Definition 4. The function $\widetilde{c} : \prod_{i=1}^n \widetilde{\Lambda}_i \rightarrow R$ defined by

$$(u_{k_1}^1, \dots, u_{k_n}^n) \mapsto \widetilde{c}(u_{k_1}^1, \dots, u_{k_n}^n) = \frac{\#\{(u_{k_1}^1, \dots, u_{k_n}^n) \in \mathcal{T}\}}{N}, k_i = 0, \dots, m_i$$

will be called sample derivated copula distribution.

Remark 2. The maximum number of different vectors what the above defined sample derivated copula can take on equals to $\prod_{i=1}^n m_i$.

Definition 5. The function $\tilde{C}_n^{\mathcal{P}} : \prod_{i=1}^n \tilde{\Lambda}_i \subset [0; 1]^n \rightarrow [0; 1]$ defined by

$$\begin{aligned} (u_{k_1}^1, \dots, u_{k_n}^n) &\mapsto \tilde{C}_n^{\mathcal{P}}(u_{k_1}^1, \dots, u_{k_n}^n) = \\ &= \frac{\#\{(u_{k_1}^1, \dots, u_{k_n}^n) \in \mathcal{T} \mid u_1 \leq u_{k_1}^1, \dots, u_n \leq u_{k_n}^n\}}{N} \end{aligned}$$

will be called sample derivated copula.

Throughout the paper we use the notation \tilde{C}_n instead of $\tilde{C}_n^{\mathcal{P}}$.

Theorem 2. *The sample derivated copula is a copula.*

PROOF. 1) It is evident that \tilde{C}_n is increasing in its each component.

2) If exists s such that $u_{k_s}^s = 0$ then $\tilde{C}_n(u_{k_1}^1, \dots, u_{k_{s-1}}^{s-1}, 0, u_{k_{s+1}}^{s+1}, \dots, u_{k_n}^n) = 0$. This follows directly from the definition. The sample do not contain any vector with a negative coordinate.

3) If for all $s \neq l$ we have $u_{k_s}^s = 1$ then

$$\begin{aligned} \tilde{C}_n(1, \dots, 1, u_{k_l}^l, 1, \dots, 1) &= \\ &= \frac{\#\{(u_1, \dots, u_s, \dots, u_n) \in \mathcal{T} \mid u_s \leq 1, \forall s \neq l \text{ and } u_l \leq u_{k_l}^l\}}{N} = \\ &= \frac{1}{N} \sum_0^{k_l-1} n_l = \frac{n_l \cdot k_l}{N} = \frac{k_l}{m_l} = u_{k_l}^l \end{aligned}$$

4) \tilde{C}_n is n -increasing as it is a cumulative probability distribution function.

Remark 3. The sample derivated copula differs, from the empirical copula [5] and the discrete copula [10], [11]. One of the differences is that the cardinal of $\tilde{\Lambda}_i$ is not necessary the same for all $i = 1, \dots, n$. Another difference is that a marginal variable can take the same value in more than one vector (since $m_i < N$).

Theorem 3. *The sample derivated copula has the following consistency property. If all variables $u_{k_s}^s = 1$ for $s \in V \setminus \{l_1, \dots, l_q\}$ then*

$$\tilde{C}_n(u_{k_1}^1, \dots, u_{k_n}^n) = \tilde{C}_q(u_{l_1}^1, \dots, u_{l_q}^q).$$

PROOF.

$$\begin{aligned} \tilde{C}_n(u_{k_1}^1, \dots, u_{k_n}^n) &= \\ &= \frac{\#\{(u_1, \dots, u_n) \in \mathcal{T} \mid u_s \leq 1, s \in V \setminus \{l_1, \dots, l_q\}, u_{l_i} \leq u_{l_i}^i, i = 1 \dots q\}}{N} = \\ &= \frac{\#\{(u_1, \dots, u_n) \in \mathcal{T} \mid u_{l_i} \leq u_{l_i}^i, i = 1 \dots q\}}{N} = \\ &= \tilde{C}_q(u_{l_1}^1, \dots, u_{l_q}^q). \end{aligned}$$

Remark 4. In general copulas do not fulfill the consistency property..

Remark 5. This theorem assures the consistency statements that we need when constructing junction tree like copulas.

At the end of this part we convince the reader from an information theoretical point of view why should one use the uniform partition and the sample derivated copula.

In the following we suppose that each $\Lambda_i, i = 1, \dots, n$ is partitioned in the same number of $m_i, i = 1, \dots, n$ intervals as in the previous case. We denote now the partitioning points by $y_j^{p_i}, j = 0, 1, \dots, m_i; i = 1, \dots, n$. This partition is arbitrary (for example equidistant) which has not the property that each interval contains the same number of sample elements. The partition of Λ_i is given by $\{y_j^{p_i} \mid j \in \{0, 1, \dots, m_i\}\}$ and is denoted by \mathcal{P}'_i . We denote by \mathcal{P}' the set of partitions $\mathcal{P}'_1, \dots, \mathcal{P}'_n$.

We denote the number of values of the variable X_i belonging to $(y_j^{p_i}; y_{j+1}^{p_i}] \cap \Lambda_i$ by k_j^i .

Let \tilde{Y}_i be the categorical random variable associated to X_i :

$$P(\tilde{Y}_i \in (y_j^{p_i}; y_{j+1}^{p_i}]) = \frac{k_j^i}{N}, j = 0, 1, \dots, m_i,$$

where

$$\sum_{j=1}^{m_i} \frac{k_j^i}{N} = 1.$$

The entropy of X_i determined by the partition \mathcal{P}'_i is:

$$H_{\mathcal{P}'_i}(X_i) = H(\tilde{Y}_i) = - \sum_{j=1}^{m_i} \frac{k_j^i}{N} \log \frac{k_j^i}{N}, \quad i = 1 \dots n.$$

It can be seen that the entropy $H_{\mathcal{P}'_i}(X_i)$ depends on the number of intervals m_i and on k_j^i .

We introduce the following notation:

$$q_{j_1, \dots, j_n}^{X_1, \dots, X_n} = P(X_1 \in (y_{j_1-1}^{p_1}; y_{j_1}^{p_1}], X_2 \in (y_{j_2-1}^{p_2}; y_{j_2}^{p_2}], \dots, X_n \in (y_{j_n-1}^{p_n}; y_{j_n}^{p_n}]),$$

where $j_i = 1, \dots, m_i$, $i = 1, \dots, n$.

The joint probability distribution determined by the partition \mathcal{P}' has the joint entropy:

$$H_{\mathcal{P}'}(X_1, \dots, X_n) = - \sum_{j_1=1}^{m_1} \dots \sum_{j_n=1}^{m_n} q_{j_1, \dots, j_n}^{X_1, \dots, X_n} \log_2 q_{j_1, \dots, j_n}^{X_1, \dots, X_n}.$$

The information content of the joint probability distribution determined by the partition \mathcal{P}' is:

$$\begin{aligned} I_{\mathcal{P}'}(X_1, \dots, X_n) &= \sum_{i=1}^n H_{\mathcal{P}'_i}(X_i) - H_{\mathcal{P}'}(X_1, \dots, X_n) = \\ &= - \sum_{i=1}^n \sum_{j=1}^{m_i} \frac{k_j^i}{N} \log \frac{k_j^i}{m_i} + \sum_{j_1=1}^{m_1} \dots \sum_{j_n=1}^{m_n} q_{j_1, \dots, j_n}^{X_1, \dots, X_n} \log_2 q_{j_1, \dots, j_n}^{X_1, \dots, X_n}. \end{aligned} \quad (7)$$

Remark 6. In this case the information content depends on the number of intervals m_i , the number of values in each one dimensional interval k_j^i , and the probabilities of belonging to the n -dimensional intervals.

If we regard the uniform partition \mathcal{P} then the entropy of X_i is:

$$H_{\mathcal{P}_i}(X_i) = H(\tilde{X}_i) = - \sum_{j=1}^{m_i} \frac{1}{m_i} \log \frac{1}{m_i} = \log m_i, \quad i = 1, \dots, n.$$

The entropy $H_{\mathcal{P}_i}(X_i)$ depends just on the number of intervals m_i .

We express now the probability

$$\begin{aligned} p_{j_1, \dots, j_n}^{X_1, \dots, X_n} &= P(X_1 \in (x_{j_1-1}^{p_1}; x_{j_1}^{p_1}], \dots, X_n \in (x_{j_n-1}^{p_n}; x_{j_n}^{p_n}]) = \\ &P(\tilde{U}_1 = u_{j_1}, \dots, \tilde{U}_n = u_{j_n}) = \tilde{c}(u_{j_1}, \dots, u_{j_n}) \end{aligned}$$

The joint probability entropy associated to the partition is:

$$\begin{aligned} H_{\mathcal{P}}(X_1, \dots, X_n) &= - \sum_{j_1=1}^{m_1} \dots \sum_{j_n=1}^{m_n} p_{j_1, \dots, j_n}^{X_1, \dots, X_n} \log_2 p_{j_1, \dots, j_n}^{X_1, \dots, X_n} = \\ &= - \sum_{j_1=1}^{m_1} \dots \sum_{j_n=1}^{m_n} \tilde{c}(u_{j_1}, \dots, u_{j_n}) \log_2 \tilde{c}(u_{j_1}, \dots, u_{j_n}). \end{aligned}$$

The information content determined by the partition \mathcal{P} is:

$$\begin{aligned} I_{\mathcal{P}}(X_1, \dots, X_n) &= \sum_{i=1}^n H_{\mathcal{P}_i}(X_i) - H_{\mathcal{P}}(X_1, \dots, X_n) = \\ &= \sum_{i=1}^n \log m_i + \sum_{j_1=1}^{m_1} \dots \sum_{j_n=1}^{m_n} \tilde{c}(u_{j_1}, \dots, u_{j_n}) \log_2 \tilde{c}(u_{j_1}, \dots, u_{j_n}) \end{aligned} \quad (8)$$

Remark 7. If we suppose that for all $i = 1, \dots, n$ the number of intervals m_i is the same for the two discussed cases then comparing formulas (7) and (8) we can see that in the case of partition \mathcal{P} the information content does not depend on the first sum of formula (8) but only on the sample derived copula.

4. The junction tree approach applied to the sample derived copula.

We introduced the sample derived copula as a discrete probability distribution with uniform marginals. We proved for this special copula in Theorem 3 that the consistency properties are fulfilled.

Let $V = \{1, \dots, n\}$ be again a set of vertices. Let be defined an acyclic hypergraph over V . We denote by Γ and \mathcal{S} the set of clusters and separators of the hypergraph which determine a junction tree \mathcal{J} . The marginal probability distributions associated to the clusters $K = \{i_1, \dots, i_t\} \in \Gamma$ are denoted by $\tilde{c}_K(\tilde{\mathbf{U}}_K) = \tilde{c}_{i_1, \dots, i_t}(\tilde{U}_{i_1}, \dots, \tilde{U}_{i_t})$. The marginal probability distributions associated to the separators are denoted in the same way by $\tilde{c}_S(\tilde{\mathbf{U}}_S)$. The joint discrete copula is shortly denoted by $\tilde{c}(\tilde{\mathbf{U}})$ and the univariate marginals by $\tilde{c}_i(\tilde{U}_i), i = 1, \dots, n$.

In this section we are going to use the following popular notation:

$$\sum_{\mathbf{u}} f(\tilde{\mathbf{U}}) = \sum_{i_1=1}^{m_1} \dots \sum_{i_n=1}^{m_n} f(\tilde{U}_1 = u_{i_1}^1, \dots, \tilde{U}_n = u_{i_n}^n),$$

where $u_{i_k}^k, i_k = 1, \dots, m_k$ are the possible values of the random variable $\tilde{U}_k, k = 1, \dots, n$ and f is an arbitrary n -dimensional function. This simplified notation is used for the products, too.

Definition 6. The junction tree distribution given by

$$\tilde{c}_{\mathcal{J}}(\tilde{\mathbf{U}}) = \frac{\prod_{K \in \Gamma} \tilde{c}_K(\tilde{\mathbf{U}}_K)}{\prod_{S \in \mathcal{S}} [\tilde{c}_S(\tilde{\mathbf{U}}_S)]^{v_S-1}}, \quad (9)$$

where v_S is the number of clusters connected by the separator S , is called copula junction tree distribution, or shortly junction tree copula.

The problem is finding the junction tree copula which fits to the sample derived copula. The goodness of fitting will be quantified by the Kullback-Leibler divergence [4].

Theorem 4. *The Kullback-Leibler divergence between the approximation (9) and the sample derived copula $\tilde{c}(\tilde{\mathbf{U}})$ is given by the formula:*

$$\begin{aligned} KL(\tilde{c}_{\mathcal{J}}(\tilde{\mathbf{U}}), \tilde{c}(\tilde{\mathbf{U}})) &= \\ &= -H(\tilde{\mathbf{U}}) - \left[\sum_{K \in \Gamma} I(\tilde{\mathbf{U}}_K) - \sum_{S \in \mathcal{S}} (v_S - 1) I(\tilde{\mathbf{U}}_S) \right] + \sum_{i=1}^n \log_2 m_i. \end{aligned}$$

PROOF.

$$\begin{aligned}
& KL\left(\tilde{c}_{\mathcal{J}}\left(\tilde{\mathbf{U}}\right), \tilde{c}\left(\tilde{\mathbf{U}}\right)\right) = \\
& = \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \frac{\tilde{c}\left(\tilde{\mathbf{U}}\right)}{\tilde{c}_{\Gamma}\left(\tilde{\mathbf{U}}\right)} = \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \tilde{c}\left(\tilde{\mathbf{U}}\right) - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \tilde{c}_{\Gamma}\left(\tilde{\mathbf{U}}\right) = \\
& = -H\left(\tilde{\mathbf{U}}\right) - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \frac{\prod_{K \in \Gamma} \tilde{c}_K\left(\tilde{\mathbf{U}}_K\right)}{\prod_{S \in \mathcal{S}}\left[\tilde{c}_S\left(\tilde{\mathbf{U}}_S\right)\right]^{v_S-1}} = -H\left(\tilde{\mathbf{U}}\right) - \\
& - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \left[\log_2 \prod_{K \in \Gamma} \tilde{c}_K\left(\tilde{\mathbf{U}}_K\right) - \log_2 \prod_{S \in \mathcal{S}}\left[\tilde{c}_S\left(\tilde{\mathbf{U}}_S\right)\right]^{v_S-1}\right] = \\
& = -H\left(\tilde{\mathbf{U}}\right) - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \prod_{K \in \Gamma} \tilde{c}_K\left(\tilde{\mathbf{U}}_K\right) + \\
& + \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \prod_{S \in \mathcal{S}}\left[\tilde{c}_S\left(\tilde{\mathbf{U}}_S\right)\right]^{v_S-1}.
\end{aligned}$$

We add and subtract the sum:

$$\sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \prod_{K \in \Gamma} \prod_{i \in K} \tilde{c}_i\left(\tilde{U}_i\right). \quad (10)$$

It follows from the definition of the junction tree that $\bigcup_{K \in \Gamma} K = V$, and each variable belongs once more in the clusters as in the separators. So by adding and subtracting (10) we obtain the following:

$$\begin{aligned}
& KL\left(\tilde{c}_{\mathcal{J}}\left(\tilde{\mathbf{U}}\right), \tilde{c}\left(\tilde{\mathbf{U}}\right)\right) = -H\left(\tilde{\mathbf{U}}\right) - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \frac{\prod_{K \in \Gamma} \tilde{c}_K\left(\tilde{\mathbf{U}}_K\right)}{\prod_{K \in \Gamma} \prod_{i \in K} \tilde{c}_i\left(\tilde{U}_i\right)} + \\
& + \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \frac{\prod_{S \in \mathcal{S}}\left[\tilde{c}_S\left(\tilde{\mathbf{U}}_S\right)\right]^{v_S-1}}{\prod_{S \in \mathcal{S}}\left[\prod_{i \in S} \tilde{c}_i\left(\tilde{U}_i\right)\right]^{v_S-1}} - \sum_{\mathbf{u}} \tilde{c}\left(\tilde{\mathbf{U}}\right) \log_2 \prod_{i=1}^n \tilde{c}_i\left(\tilde{U}_i\right) =
\end{aligned}$$

$$\begin{aligned}
&= -H(\tilde{\mathbf{U}}) - \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{K \in \Gamma} \log_2 \frac{\tilde{c}_K(\tilde{\mathbf{U}}_K)}{\prod_{i \in K} \tilde{c}_i(\tilde{U}_i)} + \\
&+ \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{S \in \mathcal{S}} \log_2 \frac{[\tilde{c}_S(\tilde{\mathbf{U}}_S)]^{v_S-1}}{\left[\prod_{i \in S} \tilde{c}_i(\tilde{U}_i) \right]^{v_S-1}} - \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{i=1}^n \log_2 \tilde{c}_i(\tilde{U}_i).
\end{aligned}$$

Since the sample derived copula has the property that all $\tilde{c}_K(\tilde{\mathbf{U}}_K)$, $\tilde{c}_S(\tilde{\mathbf{U}}_S)$, $\tilde{c}_i(\tilde{U}_i)$ are consistent marginals of $\tilde{c}(\tilde{\mathbf{U}})$ (see Theorem 3) we have the following relations:

$$\begin{aligned}
&\bullet \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{K \in \Gamma} \log_2 \frac{\tilde{c}_K(\tilde{\mathbf{U}}_K)}{\prod_{i \in K} \tilde{c}_i(\tilde{U}_i)} = \\
&= \sum_{K \in \Gamma} \sum_{\mathbf{u}_K} \tilde{c}_K(\tilde{\mathbf{U}}_K) \log_2 \frac{\tilde{c}_K(\tilde{\mathbf{U}}_K)}{\prod_{i \in K} \tilde{c}_i(\tilde{U}_i)} = \sum_{K \in \Gamma} I(\tilde{\mathbf{U}}_K); \\
&\bullet \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{S \in \mathcal{S}} \log_2 \frac{[\tilde{c}_S(\tilde{\mathbf{U}}_S)]^{v_S-1}}{\left[\prod_{i \in S} \tilde{c}_i(\tilde{U}_i) \right]^{v_S-1}} = \\
&= \sum_{S \in \mathcal{S}} \sum_{\mathbf{u}_S} (v_S - 1) \tilde{c}_S(\tilde{\mathbf{U}}_S) \log_2 \frac{\tilde{c}_S(\tilde{\mathbf{U}}_S)}{\prod_{i \in S} \tilde{c}_i(\tilde{U}_i)} = \sum_{S \in \mathcal{S}} (v_S - 1) I(\tilde{\mathbf{U}}_S); \\
&\bullet - \sum_{\mathbf{u}} \tilde{c}(\tilde{\mathbf{U}}) \sum_{i=1}^n \log_2 \tilde{c}_i(\tilde{U}_i) = \sum_{i=1}^n H(\tilde{U}_i) = \sum_{i=1}^n \log_2 m_i;
\end{aligned}$$

Here $I(\tilde{\mathbf{U}}_K)$, $I(\tilde{\mathbf{U}}_S)$ are the information content of the probability distribution of the marginals $\tilde{c}_K(\tilde{\mathbf{U}}_K)$ and $\tilde{c}_S(\tilde{\mathbf{U}}_S)$ (see [4]).

By the substitution of these assertions we obtain:

$$\begin{aligned}
KL(\tilde{c}_{\mathcal{J}}(\tilde{\mathbf{U}}), \tilde{c}(\tilde{\mathbf{U}})) &= \\
&= -H(\tilde{\mathbf{U}}) - \left[\sum_{K \in \Gamma} I(\tilde{\mathbf{U}}_K) - \sum_{S \in \mathcal{S}} (v_S - 1) I(\tilde{\mathbf{U}}_S) \right] + \sum_{i=1}^n \log_2 m_i.
\end{aligned}$$

Remark 8. The difference $\sum_{i=1}^n \log_2 m_i - H(\tilde{\mathbf{U}})$ does not depend on the junction tree structure.

Definition 7. The difference

$$\sum_{K \in \Gamma} I(\tilde{\mathbf{U}}_K) - \sum_{S \in \mathcal{S}} (v_S - 1) I(\tilde{\mathbf{U}}_S)$$

is called the weight of the junction tree copula.

It is easy to see that in order to find a better approximation using junction trees, we have to find the junction tree having the largest weight.

Finding the best fitting k -width junction tree, (the largest cluster contains k elements) for $k > 2$ is an NP-hard problem. For $k = 2$ the problem is similar to the Chow-Liu approximation [2]. In this case it is possible to find the best fitting second order junction tree by Kruskal' or Prim' algorithm.

For $k \geq 3$ it can be successfully used a heuristic approach introduced by the authors in [7] and [14]. The idea is the fitting of a special kind of junction tree, called t -cherry junction tree.

5. Conclusions and possible applications

One of the advantages of the junction tree copula is that it reveals some of the conditional independences between the variables involved. This kind of dependence structure is not exploited by the copula function. Another advantage of the method is that a multivariate copula can be decomposed into some lower dimensional sample derived copulas.

The sample derived copula approach is useful in cases when nothing else is known about the probability distribution but an iid sample. If the uniform partition is applied the whole information content depends on the sample derived copula.

The copula junction tree can be used in feature selection which is a key-question in many fields as finance, medicine and biostatistics.

We got very good numerical results in pattern recognition (see [15]). First we applied the uniform partition to discretize continuous random variables then constructed the t -cherry junction tree approximation. In this way we found the informative features and so reduced the dimension of the classifier.

References

- [1] R.S. Calsaverini, R. Vicente, An information theoretic approach to statistical dependence: Copula information, arXiv:0911.4207v1, 2009.
- [2] C.K. Chow, C.N. Liu, Approximating discrete probability distribution with dependence tree, IEEE Transactions on Information Theory 14 (1968) 462–467.
- [3] Gh. Constantin, I. Istratescu, Elements of Probabilistic Analysis, Kluwer Academic Publisher, Dordrecht, 1989.
- [4] T.M. Cover, J.A. Thomas, Elements of Information Theory, Wiley Interscience, New York, 1991.
- [5] P. Deheuvels, La fonction de dépendance empirique et ses propriétés – Un test non paramétrique d’indépendance, Académie Royale de Belgique – Bulletin de la Classe des Sciences – 5e Série 65 (1979) 274–292.
- [6] S. Kirshner, Learning with tree-average densities and distributions, Advances in Neural Information Processing Systems (NIPS), 2007.
- [7] E. Kovács, T. Szántai, On the approximation of discrete multivariate probability distribution using the new concept of t -cherry junction tree, Lecture Notes in Economics and Mathematical Systems, 633, Proceedings of the IFIP/IIASA/GAMM Workshop on Coping with Uncertainty, Robust Solutions, 2008, IIASA, Laxenburg, 39–56.
- [8] S.L. Lauritzen, T. Speed, K. Vijayan, Decomposable graphs and hypergraphs, J. Aust. Math. Soc. A 36 (1984) 12–29.
- [9] S.L. Lauritzen, D.J. Spiegelhalter, Local computations with probabilities on graphical structures and their application to expert systems, Journal of the Royal Statistical Society, Ser. B 50 (1988) No. 2 157–224.

- [10] G. Mayor, J. Suner, J. Torrens, Copula-like operations on finite settings, *IEEE Transactions on fuzzy systems* 13 (2005) No. 4 468–477.
- [11] R. Mesiar, Discrete copulas – what they are, In: *Joint EUSFLAT-LFA 2005, Conference Proceedings* (E. Montseny and P. Sobrevilla, eds.) Universitat Politècnica de Catalunya, Barcelona, 2005, 927-930.
- [12] J. Ma, Z. Sun, Mutual information is copula entropy, arXiv: 0808.0845v1, 2008.
- [13] R. Nelsen, *An Introduction to Copulas*, Springer, New York, 1999.
- [14] Szántai, T. and E. Kovács, Hypergraphs as a mean of discovering the dependence structure of a discrete multivariate probability distribution, *Proc. Conference APplied mathematical programming and MODelling (APMOD), 2008*, Bratislava, 27-31 May 2008, *Annals of Operations Research*, to appear.
- [15] T. Szántai, E. Kovács, Application of t -cherry junction trees in pattern recognition, *Broad Research in Artificial Intelligence and Neuroscience (BRAIN), Special Issue on Complexity in Sciences and Artificial Intelligence*, Eds. B. Iantovics, D. Radoiu, M. Marusteri and M. Dehmer, 2010, 40–45.
- [16] R.E. Tarjan, M. Yannakakis, Simple linear time algorithms to test chordality of graphs, test acyclicity of hypergraphs, and selectively reduce acyclic hypergraphs, *SIAM J. Comput.* 13 (1984) 566–579.
- [17] N.N. Vorob’ev, Consistent families of measures and their extensions, *Theory of Probability and Applications* 7 (1962) 147–163.
- [18] N.N. Vorob’ev, Markov measures and Markov extensions, *Theory of Probability and Applications* 8 (1963) 420–429.
- [19] J. Yang, Y. Qi, Q.R. Wang, A class of multivariate copulas with bivariate Fréchet marginal copulas, *Insurance: Mathematics and Economics* 45 (2009) 139–147.