

Modeling bivariate distributions with triangular fuzzy data and its application in hydrological studies: A copula-based approach

P. Khalilpour ¹, A. Parchami ² and R. Pourmoussa ³

^{1,2,3}Department of Statistics, Faculty of Mathematics and Computer, Shahid Bahonar University of Kerman, Kerman, Iran

parisa_khalilpour@math.uk.ac.ir, parchami@uk.ac.ir, pourm@uk.ac.ir

Abstract

Fuzzy data analysis presents significant computational challenges due to its inherent ambiguity and uncertainty. Traditional statistical methods do not have the capability to effectively capture and model the uncertainty in fuzzy observations. A novel approach is proposed in this paper to model unknown bivariate densities between variables with fuzzy observations and incorporating the dependency. By employing this copula-based approach, we have effectively managed the computational complexity associated with the analysis of fuzzy data. The proposed approach has been applied to model groundwater aquifers distribution.

Keywords: Bivariate density estimation, fuzzy observations, *FGM* copula, *AMH* copula, Gaussian copula.

1 Introduction

Statistical modeling and distribution estimation are fundamental to data analysis, enabling a deeper understanding of data structures, uncovering hidden patterns, and making more accurate predictions. By identifying appropriate statistical distributions, essential characteristics such as mean and variance can be analyzed, providing a solid foundation for selecting suitable analytical methods. In the multivariate case, statistical modeling becomes even more critical, as many real-world phenomena are influenced by multiple variables simultaneously. Multivariate models allow us to capture intricate relationships between variables and assess their mutual impact, leading to more comprehensive insights. These models are particularly valuable in fields such as economics, medicine, social sciences, and engineering, where precise decision-making relies on understanding the interconnectedness of multiple factors [16, 18]. Data modeling becomes even more important when the data are imprecise and fuzzy rather than crisp. Traditional methods often struggle with uncertainty and imprecision, making it difficult to extract valuable insights. Given that some data is inherently ambiguous or imprecise, modeling in this context becomes crucial and requires methods capable of managing such computational complexities [21]. Copulas play a crucial role in fuzzy data modeling, as they provide a powerful framework for analyzing dependencies in uncertain and imprecise datasets. Unlike traditional methods that often assume rigid linear relationships, copulas allow for the modeling of nonlinear and complex dependencies, enabling a more accurate representation of the underlying structure of fuzzy data. Moreover, they facilitate the combination of different marginal distributions without imposing restrictive assumptions on variable distributions.

Arefi et al. [1] employed three methods, that were extended for density estimation based on α -cuts of fuzzy random variables. Hesamian and Akbari [4] extended classical nonparametric curve-fitting methods to fuzzy random variables, covering fuzzy density estimation and nonparametric regression. Lukasik et al. [13] presented an automatic, unsupervised clustering algorithm for synthesizing fuzzy models. Khorramdel et al. [8] proposed a framework using a diffusion-based kernel density estimator with adaptive bandwidth selection to generate high-quality prediction intervals for non-stationary wind power time series. Sun and Khayatnezhad [20] used fuzzy set theory and copulas, like the

Gumbel copula, to model nonlinear dependencies between rainfall and flood characteristics. Mortuza et al. [15] used fuzzy clustering and copula models to analyze the joint probability distribution of drought duration and severity in Bangladesh. Nelsen [16] provided a comprehensive guide on bivariate copulas, while Joe [6] included a chapter on multivariate copulas in his book. He [16] explained the Farlie-Gumbel-Morgenstern (*FGM*) family by looking at the correlation between the marginals. Huang and Kotz [5] created a polynomial-type extension of the *FGM* bivariate distributions using a single parameter. They also introduced a new way to characterize bivariate copulas using Dini derivatives. Kim et al. [9] proposed a new class of bivariate copulas to measure dependence and include them in various iterated copula families. Kumar [10] demonstrated that the Pearson product-moment correlation models only model linear dependence and assume normal distributions, whereas a copula-based approach, such as the Ali-Mikhail-Haq (*AMH*) copula, provides more flexibility in modeling multivariate distributions without relying on normality or independence assumptions. Renard et al. [17] explored the application of the Gaussian copula in multivariate frequency analysis for hydrological risk assessment, highlighting its limitations and potential errors. Meyer [14] translated well-known and lesser known facts about the bivariate normal distribution into copula language, providing various expressions for the bivariate normal copula, calculating its Gini's gamma, and deriving improved bounds and approximations on its diagonal.

The motivation behind this study is to extend the applicability of statistical modeling to fuzzy data, addressing the limitations of traditional methods. By integrating copulas with fuzzy observations, we aim to develop a more robust approach that effectively captures dependencies in uncertain environments. In this paper, we have selected three copulas, *FGM*, *AMH* and Gaussian, as the basis for our analysis to effectively explore the complex relationships between variables with fuzzy observations. This not only enhances the accuracy of data analysis but also broadens the scope of statistical modeling, making it more adaptable to real-world complexities.

This paper is structured as follows. We begin our study of copulas for bivariate distributions with precise observations and review famous copulas and Sklar's theorem in Section 2. In Section 3, we discuss bivariate distributions for variables with fuzzy observations, and the value of the correlation between random variables is obtained. The evaluation of the copula-based approach for bivariate fuzzy data is presented in Section 4. An application in hydrological studies is discussed in Section 5, which shows the importance of using a copula in the proposed bivariate density estimation method. The final section provides an overall conclusion.

2 Copula preliminaries

In probability theory and statistics, a copula is essentially a multivariate cumulative distribution function (*cdf*) in which the marginal distribution of each variable is uniform on the interval $[0, 1]$. Copulas are used as a powerful tool for modeling the dependence between random variables. In high-dimensional statistical analysis, copulas play a crucial role as they link the marginal distributions of individual variables to their joint distribution. There are various parametric families of copulas, typically characterized by parameters that control the strength of dependence. In the following, some commonly used parametric copula models are introduced. The *FGM* copula is used to model weak dependence between variables, as its implementation is very simple. The *FGM* copula formula for two random variables X and Y with marginals F_X and F_Y is defined as:

$$C(F_X(x), F_Y(y)) = F_X(x)F_Y(y) (1 + \theta(1 - F_X(x))(1 - F_Y(y))), \quad (1)$$

where $\theta \in [-1, 1]$ is the correlation dependence parameter between two random variables X and Y [11, 16]. Function $C : [0, 1]^2 \rightarrow [0, 1]$ is actually the bivariate *cdf*, and to obtain the bivariate probability density function (*pdf*), we need to take the partial derivatives with respect to F_X and F_Y . Another widely used copula for dependence modeling is the *AMH*, which accommodates both positive and negative dependencies. With only one parameter controlling the strength of dependence, the *AMH* copula enables efficient computational handling. The *AMH* copula for two random variables X and Y on the domain $C : [0, 1]^2$ is defined by [11, 16]:

$$C(F_X(x), F_Y(y)) = \frac{F_X(x)F_Y(y)}{1 - \theta(1 - F_X(x))(1 - F_Y(y))}. \quad (2)$$

The Gaussian copula is another widely used tool in statistical modeling, particularly for capturing complex dependence structures in multivariate distributions. Its strength lies in its ability to represent varying degrees of correlation, accommodating both linear and non-linear dependencies. The Gaussian copula is derived from the multivariate normal distribution, making it a flexible framework for modeling bivariate behaviors of random variables. Mathematically, the Gaussian copula for two random variables X and Y on $C : [0, 1]^2 \rightarrow [0, 1]$ is given by [11, 16]:

$$C(F_X(x), F_Y(y)) = \Phi_2[\Phi^{-1}(F_X(x)), \Phi^{-1}(F_Y(y))], \quad (3)$$

where Φ_2 is the bivariate *cdf* of a standard bivariate normal distribution with correlation parameter θ and Φ^{-1} is the quantile function of the standard normal distribution.

One of the most important applications of copulas in probability theory and statistics is establishing the relationship between two or more variables. This application becomes crucial when we want to model and analyze complex dependencies between variables. According to Sklar's theorem, given the univariate marginal distributions of each variable and combining them with a copula that captures the dependence structure between the variables, one can easily obtain the bivariate *cdf*. Based on Sklar's Theorem [16], if X and Y are two random variables with a bivariate *cdf* $F_{X,Y}$ and marginal *cdfs* F_X and F_Y , respectively, then there exists a copula function C such that

$$F_{X,Y}(x,y) = C(F_X(x), F_Y(y)), \text{ for all } (x,y) \in R^2. \quad (4)$$

If F_X and F_Y are continuous, then C is unique; otherwise, it is uniquely determined on $Ran(F_X) \times Ran(F_Y)$ where $Ran(F)$ is the range of the distribution function F . Conversely, if C is a copula and F_X and F_Y are *cdfs*, then the function F defined by Eq. (4) is the bivariate *cdf* of (X, Y) with margins F_X and F_Y . Refer to references [11, 16] for the conditions and properties of copula functions.

3 Estimation of bivariate distributions by fuzzy observations

One of the fundamental and common applications of the copula and Sklar's Theorem is obtaining the bivariate *cdf*. Estimating both the univariate distributions and the bivariate distribution is a challenging statistical problem, we are going to investigate bivariate density estimation when data are fuzzy/imprecise rather than crisp. A novel approach is proposed for estimating the marginal distributions and determining the dependence structure between continuous random variables with fuzzy observations. Then, the copula and Sklar's Theorem, by establishing a connection between these marginal distributions, serve as powerful tools for modeling the dependence between random variables. This approach involves the following three steps:

Step 1 (Estimating marginal distributions): To estimate univariate marginal distributions of continuous random variables X and Y , one can utilize any possible approach. Here, a special method is proposed for this aim based on aggregation observed membership degrees based on the mean function for each point, which is a generalized fuzzy-data-based version of the kernel density estimation [7]. In other words, the univariate density estimation based on the vector of fuzzy observations $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ is defined in [7, 22] by

$$\hat{f}_{\tilde{\mathbf{x}}}(t) = \overline{\tilde{x}^*(t)} = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i^*(t), \quad (5)$$

where $\tilde{x}_i^*(t) = \frac{\tilde{x}_i(t)}{\int_{-\infty}^{+\infty} \tilde{x}_i(t) dt}$ is the normalized fuzzy number \tilde{x}_i for $i = 1, \dots, n$. Therefore, the univariate estimation method based on fuzzy data in point x is equal to the arithmetic mean over the normalized membership values at point x for fuzzy data. It must be noted that this averaging differs from the fuzzy arithmetic mean and does not rely on Zadeh's extension principle. Moreover, the estimation of *cdf* of a continuous random variable X based on fuzzy observations is equal to

$$\hat{F}_{\tilde{\mathbf{x}}}(t) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^t \tilde{x}_i^*(s) ds. \quad (6)$$

But before starting the example, we define notation $T(a_i, b_i, c_i)$ for a triangular fuzzy number corresponding to the i th observation with the membership function

$$T(a_i, b_i, c_i)(x) = \begin{cases} \frac{x-a_i}{b_i-a_i} & a_i \leq x < b_i \\ \frac{c_i-x}{c_i-b_i} & b_i \leq x < c_i \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

To better understand and clarify the proposed univariate density estimation method, see the below numerical example.

Remark 3.1. *In this study, the underlying variables of interest are fundamentally classical/crisp random variables, while their observed values for these variables are recorded as fuzzy sets. This approach is motivated by the*

presence of epistemic uncertainty, often arising from imprecise measurement tools, the subjective nature of human observation, or the inherent difficulty in defining sharp thresholds for continuous phenomena. For example, the true salinity at a depth of 20 meters of soil is a precise and certain value, the instrumental inaccuracy of a salinity sensor introduces an error that causes the salinity level to be inaccurately recorded. (See references [1] and [4] for such cases.) Consequently, representing the observed value as a fuzzy number provides a more rigorous and truthful mathematical framework for handling this imperfection in the data, allowing for a more robust analysis that explicitly accounts for measurement uncertainty.

Example 3.2. Consider six triangular fuzzy observations $\tilde{x}_1 = T(0.8, 2, 3)$, $\tilde{x}_2 = T(1, 3, 5)$, $\tilde{x}_3 = T(2, 4, 6)$, $\tilde{x}_4 = T(3, 5.5, 6)$, $\tilde{x}_5 = T(3, 3.5, 4.3)$ and $\tilde{x}_6 = T(4, 4.5, 5.5)$ for a continuous random variable X (see the first image in Figure 1). Using Eqs. (5-6), the estimation of the *PDF* and *CDF* for random variable X is depicted, respectively in the second and third images of Figure 1. For instance, according to Eq. (5), the univariate density estimation method can be obtained at point 4 via averaging over the normalized membership functions of fuzzy data $\tilde{x}_1, \dots, \tilde{x}_6$ as follows

$$\begin{aligned} \hat{f}_{\tilde{x}}(4) &= \frac{1}{6} \sum_{i=1}^6 \tilde{x}_i^*(4) \\ &= \frac{1}{6} \left[\frac{\tilde{x}_1(4)}{1.1} + \frac{\tilde{x}_2(4)}{2} + \frac{\tilde{x}_3(4)}{2} + \frac{\tilde{x}_4(4)}{1.5} + \frac{\tilde{x}_5(4)}{0.8} + \frac{\tilde{x}_6(4)}{1} \right] \\ &= \frac{0 + 0.25 + 0.5 + 0.26 + 0.57 + 0}{6} = 0.26, \end{aligned}$$

which means that the *pdf* value of the random variable X at point $x = 4$ is estimated by 0.26. Based on Eq. (6) and the following calculations, the *CDF* value at point $x = 4$ is obtained as 0.56,

$$\begin{aligned} \hat{F}_{\tilde{x}}(4) &= \frac{1}{6} \sum_{i=1}^6 \int_{-\infty}^4 \tilde{x}_i^*(s) ds \\ &= \frac{1}{6} \left[\frac{\int_{-\infty}^4 \tilde{x}_1(s) ds}{1.1} + \frac{\int_{-\infty}^4 \tilde{x}_2(s) ds}{2} + \frac{\int_{-\infty}^4 \tilde{x}_3(s) ds}{2} + \frac{\int_{-\infty}^4 \tilde{x}_4(s) ds}{1.5} + \frac{\int_{-\infty}^4 \tilde{x}_5(s) ds}{0.8} + \frac{\int_{-\infty}^4 \tilde{x}_6(s) ds}{1} \right] \\ &= \frac{1 + 0.87 + 0.5 + 0.13 + 0.91 + 0}{6} = 0.56. \end{aligned}$$

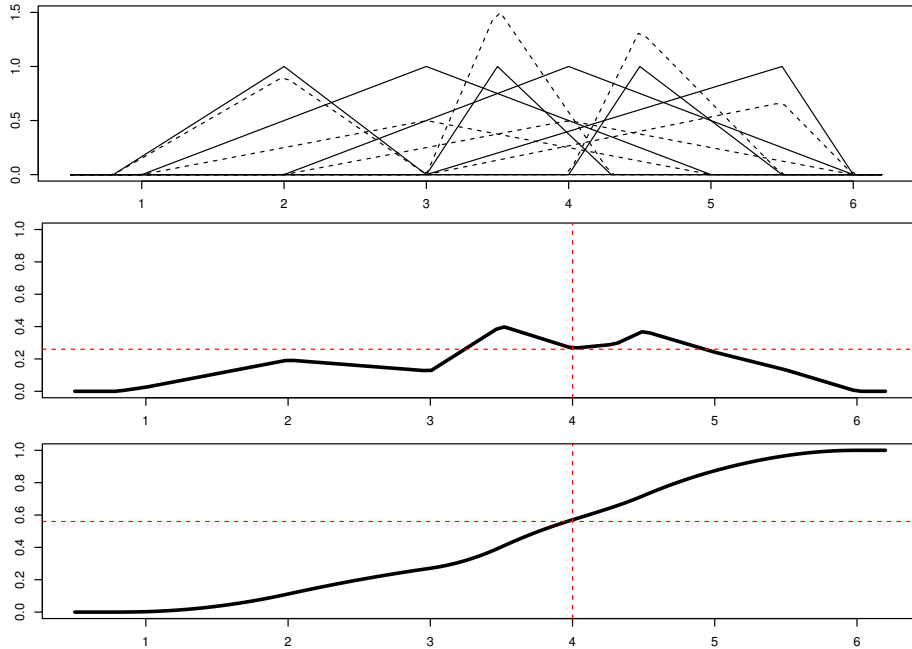


Figure 1: First image: Membership functions of fuzzy observations (continuous lines), normalized membership functions of fuzzy observations (dashed lines), second image: estimation of *pdf* and the last image is the estimation of *cdf* based on fuzzy observations in Example 3.2.

Step 2 (Modeling dependencies with copulas): Dependence structure analysis and relationships between variables are one of the important and widely studied issues in statistics. Calculating the correlation coefficient between two sets of triangular fuzzy data is a challenging statistical problem. To simplify this process, we propose a method based on Pearson's correlation coefficient (corr) in which three following correlation coefficients are calculated using parameters of triangular fuzzy numbers: (1) Correlation coefficient between the cores of the triangular fuzzy observations \tilde{x}_i 's and \tilde{y}_i 's, (2) Correlation coefficient between the lower bounds of the supports of the triangular fuzzy observations \tilde{x}_i 's and \tilde{y}_i 's, and (3) Correlation coefficient between the upper bounds of the supports of the triangular fuzzy observations \tilde{x}_i 's and \tilde{y}_i 's. Then, the average of these calculated correlation coefficients is considered as an estimate of the dependency parameter ($\hat{\theta}$) between two variables with fuzzy observations. It is worth noting that, without loss of generality, other methods can also be employed in this study to estimate the correlation between \tilde{x}_i 's and \tilde{y}_i 's in a way that ensures crisp estimation.

Step 3 (Applying Sklar's theorem): Sklar's theorem is a fundamental result that establishes a connection between bivariate *cdfs* and their marginals through a copula. In this step, it is enough to replace the estimated *cdfs* $\hat{F}_{\tilde{x}}$ and $\hat{F}_{\tilde{y}}$ from Step 1 in Eq. (6). To this aim, the following definition is presented.

Definition 3.3. *The estimation of CDF for a bivariate random variable (X, Y) based on fuzzy observations (\tilde{x}, \tilde{y}) is*

$$\hat{F}_{\tilde{x}, \tilde{y}}(x, y) = C(\hat{F}_{\tilde{x}}(x), \hat{F}_{\tilde{y}}(y)), \quad (8)$$

in which $\hat{F}_{\tilde{x}}(x) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^x \tilde{x}_i^*(t) dt$ is *cdf* estimation of random variable X based on fuzzy observations $\tilde{x}_1, \dots, \tilde{x}_n$. The corresponding bivariate *pdf* can be obtained by differentiating $\hat{F}_{\tilde{x}, \tilde{y}}(x, y)$ with respect to x and y as follows

$$\hat{f}_{\tilde{x}, \tilde{y}}(x, y) = \left(1 + \hat{\theta}(1 - 2\hat{F}_{\tilde{x}}(x))(1 - 2\hat{F}_{\tilde{y}}(y))\right) \hat{f}_{\tilde{x}}(x) \hat{f}_{\tilde{y}}(y). \quad (9)$$

The proposed steps for bivariate *cdf* estimation based on fuzzy observations is gathered in the following algorithm for computational processes.

Proposition 3.4. *The function $\hat{F}_{\tilde{x}, \tilde{y}}$ defined in Eq. (8) is a bivariate *cdf* of (X, Y) based on fuzzy observations (\tilde{x}, \tilde{y}) with the following properties:*

- 1) *Non-negativity:* $\hat{F}_{\tilde{x}, \tilde{y}}(x, y) \geq 0$ for all $(x, y) \in R^2$,
- 2) *Boundary conditions:* $\hat{F}_{\tilde{x}, \tilde{y}}(x, \infty) = \hat{F}_{\tilde{x}}(x)$, $\hat{F}_{\tilde{x}, \tilde{y}}(\infty, y) = \hat{F}_{\tilde{y}}(y)$, $\hat{F}_{\tilde{x}, \tilde{y}}(\infty, \infty) = 1$, $\hat{F}_{\tilde{x}, \tilde{y}}(-\infty, y) = 0$ and $\hat{F}_{\tilde{x}, \tilde{y}}(x, -\infty) = 0$.
- 3) *Monotonicity:* The function $\hat{F}_{\tilde{x}, \tilde{y}}$ is non-decreasing in both x and y .

Theorem 3.5. *The integration of the introduced bivariate *pdf* estimation in Eq. (9), over the unit square $[0, 1]^2$, is equal to one.*

Proof. Since $\int_0^1 \int_0^1 \hat{f}_{\tilde{x}, \tilde{y}}(x, y) dx dy = \hat{F}_{\tilde{x}, \tilde{y}}(\infty, \infty) = 1$, the theorem is proved. \square

The presented Table 1 illustrates a broad spectrum of examples where dependent fuzzy-valued random variables play a critical role in modeling uncertainty across various scientific and engineering disciplines. Each example highlights pairs of interconnected variables whose measurements and interactions cannot be precisely measured due to inherent vagueness, measurement limitations, or environmental variability. The inclusion of diverse domains such as hydrology, medicine, astrophysics, and management underscores the ubiquitous nature of fuzzy random variables and the necessity of sophisticated probabilistic approaches to accurately capture and analyze real-world phenomena. Moreover, these examples demonstrate the practical applicability of joint distributions based on fuzzy observation in addressing complex problems involving uncertainty and dependency. For instance, estimation and understanding the joint PDF for daily rainfall intensity and surface runoff is crucial for flood management. This comprehensive overview emphasizes the relevance of fuzzy data as indispensable tools for estimation, decision-making, and better understanding the different environments in multidisciplinary contexts.

Algorithm 1 Procedure of bivariate distribution estimation**Require:**

Dataset of bivariate fuzzy observations $(\tilde{x}_i, \tilde{y}_i)$ for $i = 1, \dots, n$, where $\tilde{x}_i = T(x_{1i}, x_{2i}, x_{3i})$ and $\tilde{y}_i = T(y_{1i}, y_{2i}, y_{3i})$.

Ensure:

Bivariate *cdf* $\hat{F}_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}$.

Computations:**Step 1 (Estimating marginal distributions):**

for $i = 1$ to n do

 Compute normalized fuzzy observations $\tilde{x}_i^*(t)$ and $\tilde{y}_i^*(t)$.

end for

Estimate marginal *cdfs*: $\hat{F}_{\tilde{\mathbf{x}}}(t) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^t \tilde{x}_i^*(s) ds$ and $\hat{F}_{\tilde{\mathbf{y}}}(t) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^t \tilde{y}_i^*(s) ds$.

Step 2 (Computing dependency parameter):

for $j = 1$ to 3 do

$x^j = (x_{j1}, \dots, x_{jn})$

$y^j = (y_{j1}, \dots, y_{jn})$

$r_j = \text{corr}(x^j, y^j)$

end for

Estimate the dependency parameter by $\hat{\theta} = \text{mean}(r_1, r_2, r_3)$.

Step 3 (Applying Sklar's theorem):

Note: If the estimated parameter $\hat{\theta}$ is negative, this algorithm does not support such cases. One should instead consider copulas that allow for negative dependence, such as rotated Archimedean copulas.

if $0 \leq \hat{\theta} < 0.15$ then, use the **FGM** copula by

$$\hat{F}_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(x, y) = \hat{F}_{\tilde{\mathbf{x}}}(x) \hat{F}_{\tilde{\mathbf{y}}}(y) \left(1 + 3\hat{\theta}(1 - \hat{F}_{\tilde{\mathbf{x}}}(x))(1 - \hat{F}_{\tilde{\mathbf{y}}}(y)) \right)$$

else if $0.15 \leq \hat{\theta} < 0.33$ then use **AMH** copula by

$$\hat{F}_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(x, y) = \frac{\hat{F}_{\tilde{\mathbf{x}}}(x) \hat{F}_{\tilde{\mathbf{y}}}(y)}{1 - \frac{3\hat{\theta}\sqrt{9-36\hat{\theta}}}{2}(1 - \hat{F}_{\tilde{\mathbf{x}}}(x))(1 - \hat{F}_{\tilde{\mathbf{y}}}(y))}$$

else if $0.33 \leq \hat{\theta} \leq 1$ then use **Gaussian** copula by

$$\hat{F}_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(x, y) = \Phi_2[\Phi^{-1}(\hat{F}_{\tilde{\mathbf{x}}}(x)), \Phi^{-1}(\hat{F}_{\tilde{\mathbf{y}}}(y))]$$

end if

Return $\hat{F}_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(x, y)$.

4 Evaluation of the copula-based approach for bivariate fuzzy data

To assess the accuracy of the copula-based approach in estimating the *cdf* for bivariate fuzzy data, we conducted a focused yet limited comparison with the empirical approach. This comparison evaluates the degree of agreement between the *cdf* obtained via the copula-based method and the empirical *cdf* derived directly from fuzzy observations. Such an evaluation offers valuable insights into the effectiveness of the proposed method, especially in contexts where computational simplicity and efficiency are important. The proposed method aims to simplify the use of copulas for fuzzy data; however, such simplification mustn't significantly reduce the accuracy. To validate this, the copula-based estimates are compared with empirical distributions obtained based on fuzzy data.

In reference [1], the empirical cumulative distribution function and the empirical density function are introduced for the univariate case. Inspired by these formulations, we extend the underlying concepts to the bivariate setting and construct the corresponding empirical bivariate distribution and density functions. Using the minimum t-norm and its properties, which, unlike the product t-norm, preserve minimal dependency, we present the following definition for the bivariate case.

Definition 4.1. Let's consider the bivariate random variable (X, Y) based on fuzzy observations $\tilde{\mathbf{x}} = (\tilde{x}_1, \dots, \tilde{x}_n)^T$ with membership functions $T(a_i, b_i, c_i)$ and $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_n)^T$ with membership functions $T(a'_i, b'_i, c'_i)$, respectively. The empirical PDF (\hat{f}_n) and CDF (\hat{F}_n) are defined respectively as follows:

$$\hat{f}_n(x, y) = \frac{1}{4nh_x h_y} \sum_{i=1}^n \min\{T(a_i, b_i, c_i)(x), T(a'_i, b'_i, c'_i)(y)\}, \quad (10)$$

Table 1: Illustrative cases for paired dependent fuzzy-valued variables in various scientific domains

Scientific Domain	Variable 1 (X_1)	Variable 2 (X_2)
Hydrology	Daily rainfall intensity (mm)	Surface runoff (m ³ /s)
Structural engineering	live load on bridge (tons)	Mid-span deflection (mm)
Wind energy	Wind speed (m/s)	Turbine power output (kW)
Macroeconomics	Annual inflation rate (%)	Unemployment rate (%)
Climatology	Summer temperature (°C)	Annual evapotranspiration (mm)
Dentistry	Annual tooth decay rate	Number of dental visits
Medicine	Patient drug dosage (mg)	Immune response level
Zoology	Population of a specific species	Predation level
Veterinary medicine	Animal weight (kg)	Required vaccine dose (cc)
Biology	Cellular growth rate under specific conditions	Oxygen consumption rate
Industrial engineering	Daily order quantity	Total production time (hours)
Management	Customer satisfaction level	Brand loyalty level
Agriculture	Fertilizer amount (kg/ha)	Crop yield (tons/ha)
Computer science	Server requests per Minute	System response time (ms)
Artificial intelligence	Neural network depth	Test accuracy (%)
Electrical engineering	Input voltage (V)	Device power consumption
Geology	Well depth (m)	Subsurface pressure (bar)
Nutrition science	Daily caloric intake	Monthly body weight change
Astronomy	Distance of star to earth (light years)	Apparent luminosity
Astrophysics	Star mass (solar masses)	Star surface temperature (K)
Transportation	Instantaneous traffic (vehicles/hour)	Urban travel time (minutes)
Environmental science	Air pollution level (PM2.5)	Respiratory disease rate
Civil engineering	Foundation depth (m)	Final structure settlement (mm)
Metallurgy	Carbon percentage in steel	Final alloy hardness
Microeconomics	Product price	market demand
Education	Weekly study hours of student	Final exam score
Remote sensing	Earth surface reflectance in iR Band	Soil surface moisture
Statistics	Sample size	Accuracy of mean estimation
Robotics	Motor torque (Nm)	Robot arm speed (deg/s)
Telecommunications	Input signal power (dbm)	Bit error rate (Ber)
Social sciences	Monthly household income	Life satisfaction level
Biomedical engineering	Voltage generated by heart sensor	Final ecg signal
Psychology	Stress level	Quality of night sleep
Law	Duration of case processing (months)	Final legal cost
History	Duration of historical dynasty (years)	Number of structural reforms
Meteorology	Air pressure (hpa)	Probability of rain in next 24 hours (%)
Culture and arts	Number of gallery visitors	Art sales volume
Petroleum engineering	Reservoir pressure (psi)	Daily extraction rate (barrels)
Accounting	Number of financial documents monthly	Possible audit errors
Financial management	Initial project investment	Expected net profit (million IRR)

$$\hat{F}_n(x, y) = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^x \int_{-\infty}^y \min\{T(a_i, b_i, c_i)(t), T(a'_i, b'_i, c'_i)(s)\} dt ds, \quad (11)$$

where h_x and h_y denote the smoothing parameters associated with each variable, and $T(a_i, b_i, c_i)(x)$ and $T(a'_i, b'_i, c'_i)(y)$ represent the corresponding membership functions for the fuzzy observations of each variable as presented in Eq. (7).

In reference [12], the Anderson–Darling statistic (A^2) is introduced for the univariate case, where its structure is based on the discrepancy between the empirical distribution function and the theoretical distribution function. Inspired by this formulation, we propose an A^2 for the bivariate case.

Definition 4.2. To compare \hat{F}_n with the estimated cdf $\hat{F}_{\bar{x}, \bar{y}}$, via the copula method, we use the following A^2 to evaluate

the degree of agreement between \hat{F}_n and $\hat{F}_{\tilde{x},\tilde{y}}$

$$A^2 = n \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \frac{\left(\hat{F}_n(x, y) - \hat{F}_{\tilde{x},\tilde{y}}(x, y)\right)^2}{\hat{F}_{\tilde{x},\tilde{y}}(x, y)(1 - \hat{F}_{\tilde{x},\tilde{y}}(x, y))} d\hat{F}_n(x, y), \quad (12)$$

in which $\hat{F}_{\tilde{x},\tilde{y}}$ can be calculated by Eq. (8). Smaller values of the proposed A^2 indicate closer agreement between the empirical and copula-based bivariate fuzzy distributions, and hence reflect improved estimation accuracy.

Example 4.3. To further analyze the stability and reliability of the proposed copula-based method, its convergence behavior is examined in this example. This refers to how the estimation error changes with increasing sample size (see Table 2). For the case of $n = 10$, the triangular fuzzy data are given as follows: $T(1.42, 1.92, 2.42)$, $T(2.67, 3.17, 3.67)$, $T(0.67, 1.17, 1.67)$, $T(1.66, 2.16, 2.66)$, $T(0.91, 1.41, 1.91)$, $T(2.01, 2.50, 3.02)$, $T(1.58, 2.08, 2.58)$, $T(1.06, 1.56, 2.06)$, $T(2.23, 2.73, 3.23)$, $T(1.84, 2.34, 2.84)$, the A^2 value for these observations is approximately 0.118. Specifically, convergence means that the copula-based *cdf* tends to better approximate the empirical *cdf* as the number of fuzzy observations increases. In this example, synthetic datasets with triangular membership functions and varying sample sizes are considered. The empirical distribution is computed as defined in Eq. (8), and the copula-based *cdf* is estimated using a triangular kernel with smoothing parameters $h_x = h_y = 0.4$. The A^2 is calculated for each sample. According to the

Table 2: Values of A^2 statistic across different sample sizes

Sample size (n)	10	30	50	100	200	500	2000	5000
A^2 value	0.118	0.072	0.045	0.026	0.014	0.011	0.006	0.003

Table 2, the A^2 statistic is not only a numerical measure of goodness-of-fit but also directly interpretable [12]: the closer this value is to zero, the better the estimated density function fits the fuzzy empirical data. For larger sample sizes, a decrease in the A^2 value indicates the stabilization of the model estimation and demonstrates that the proposed method is statistically reliable. In particular, for $n = 5000$ the statistic shows desirable convergence to zero (See Figure 2). The decreasing trend of A^2 values clearly demonstrates that the discrepancy between the copula-based and empirical *cdfs* diminishes with larger sample sizes. This confirms the statistical consistency and convergence behavior of the proposed copula-based method. Moreover, the analysis strengthens the conclusion that the copula-based method not only provides a valid estimation of bivariate fuzzy distributions but also exhibits strong convergence properties as the dataset grows. As a result, the proposed method demonstrates both accuracy and reliability in practical applications.

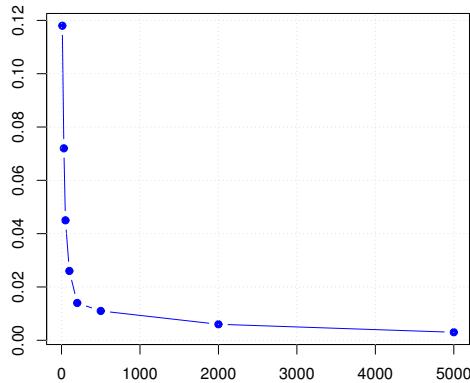


Figure 2: Convergence of A^2 statistic with increasing sample size.

5 Application in hydrological studies

Groundwater studies, especially in areas with many drilled wells, are very important. Finding the location and the estimation of underground water volume is essential for managing water resources effectively. Fifteen deep wells were randomly selected in an area with geographic coordinates between 56.7 to 57.1 longitude and 30.1 to 30.5 latitude, related to *Kerman* city in *Iran*. Suppose that the geographical coordinates and depths of wells are like the information in Table 3.

Table 3: Geographical coordinates (longitude and latitude) and precise depth (in terms of *km*) for fifteen wells in *Kerman*.

<i>i</i>	Longitude (<i>x</i>)	Latitude (<i>y</i>)	Depth (<i>o</i>)
1	56.71	30.14	0.346
2	56.73	30.21	0.340
3	56.74	30.31	0.319
4	56.75	30.41	0.329
5	56.76	30.13	0.252
6	56.83	30.48	0.297
7	56.85	30.45	0.325
8	56.95	30.43	0.271
9	56.96	30.16	0.218
10	56.97	30.48	0.273
11	56.98	30.13	0.264
12	57.05	30.30	0.291
13	57.07	30.38	0.291
14	57.09	30.33	0.286
15	57.10	30.25	0.265

One important feature of groundwater aquifers is that water is not only found exactly at the well location but also spreads around it. This means water resources can extend to areas beyond the well itself. However, the exact location of this water is uncertain and hard to pinpoint. This uncertainty is caused by factors like the geological structure of the area, how permeable the soil layers are, and local changes in groundwater flow. In continuation of this investigation, we are going to model the value of the groundwater aquifer by triangular fuzzy numbers which provides a helpful way to better understand how these groundwater resources are distributed and spread under the ground.

Non-precise coordinates of groundwater aquifers: For each sample point in the study area, the approximate coordinates of groundwater aquifers are determined/modeled using precise data of corresponding wells. The geographical coordinates longitude and latitude of each well along with its depth are utilized as reference points for this modeling. For instance, the coordinates and depth of 5th well are precisely recorded in Table 3. According to the data from 5th well, we know exactly that groundwater precisely exists at a depth of the 0.252 *km* at longitude 56.76 and latitude 30.13. It is obvious that there is also a possibility of groundwater presence around this well, and the farther we move away from the well's location (longitude and latitude), the lower this possibility becomes. Accordingly, groundwater information is modeled using fuzzy numbers in this study. In this fuzzy data modeling of the degree of possibility of groundwater presence at the longitude and latitude of each well is considered equal to one (as recorded in the Table 4), since we know that water definitely exists at this location. Therefore, the longitude and latitude of each well are regarded as the core of triangular fuzzy numbers which represent the non-precise/fuzzy location of groundwater. On the other hand, a positive coefficient of each well's depth is considered as the spread of the corresponding triangular fuzzy numbers. In other words, based on the precise data of the *i*-th well with geological coordinate (x_i, y_i) and depth o_i , the fuzzy longitude is defined as $T(x_i - ko_i, x_i, x_i + ko_i)$ and the fuzzy latitude as $T(y_i - ko_i, y_i, y_i + ko_i)$, representing the approximate groundwater coordinates in the aquifer, (see Table 4).

Selecting an appropriate value for *k* is one of the key steps in this investigation, for example for $k = 6$ which is determined by an expert, extracted data for 5th well can be modeled by fuzzy longitude $T(55.24, 56.76, 58.27)$ and fuzzy latitude $T(28.61, 30.13, 31.64)$ which essentially represents the approximate location of groundwater in the aquifer surrounding 5th well.

Univariate density estimation for marginal distributions: Now it's time to estimate the marginal *cdfs* $\hat{F}_{\bar{x}}$ and $\hat{F}_{\bar{y}}$ for the geographical coordinates of groundwater aquifers based on fuzzy observations using Eq. (6). For this purpose, the following univariate *cdfs* can be calculated for arbitrary values *x* and *y*

$$\hat{F}_{\bar{x}}(x) = \frac{1}{15} \sum_{i=1}^{15} \int_{-\infty}^x \tilde{x}_i^*(s) ds = \frac{1}{15} \left[\frac{\int_{-\infty}^x \tilde{x}_1(s) ds}{2.07} + \dots + \frac{\int_{-\infty}^x \tilde{x}_{15}(s) ds}{1.59} \right], \quad (13)$$

$$\hat{F}_{\bar{y}}(y) = \frac{1}{15} \sum_{i=1}^{15} \int_{-\infty}^y \tilde{y}_i^*(s) ds = \frac{1}{15} \left[\frac{\int_{-\infty}^y \tilde{y}_1(s) ds}{2.07} + \dots + \frac{\int_{-\infty}^y \tilde{y}_{15}(s) ds}{1.72} \right]. \quad (14)$$

Table 4: Non-precise geographical coordinates for the place of groundwater aquifers based on $k = 6$.

i	Fuzzy longitude (\tilde{x}_i)	Fuzzy latitude (\tilde{y}_i)
1	$T(54.63, 56.71, 58.78)$	$T(28.06, 30.14, 32.21)$
2	$T(54.68, 56.73, 58.77)$	$T(28.16, 30.21, 32.25)$
3	$T(54.82, 56.74, 58.65)$	$T(28.39, 30.31, 32.22)$
4	$T(54.77, 56.75, 58.72)$	$T(28.43, 30.41, 32.38)$
5	$T(55.24, 56.76, 58.27)$	$T(28.61, 30.13, 31.64)$
6	$T(55.04, 56.83, 58.61)$	$T(28.69, 30.48, 32.26)$
7	$T(54.89, 56.85, 58.80)$	$T(28.49, 30.45, 32.40)$
8	$T(55.32, 56.95, 58.57)$	$T(28.80, 30.43, 32.05)$
9	$T(55.26, 56.96, 58.65)$	$T(28.46, 30.16, 31.85)$
10	$T(55.33, 56.97, 58.60)$	$T(28.84, 30.48, 32.11)$
11	$T(55.39, 56.98, 58.56)$	$T(28.54, 30.13, 31.71)$
12	$T(55.30, 57.05, 58.79)$	$T(28.55, 30.30, 32.04)$
13	$T(55.32, 57.07, 58.81)$	$T(28.63, 30.38, 32.12)$
14	$T(55.36, 57.09, 58.81)$	$T(28.65, 30.25, 31.84)$
15	$T(55.50, 57.10, 58.69)$	$T(28.60, 30.33, 32.05)$

Bivariate density estimation of groundwater aquifers: In this step, we first compute the dependency parameter between the fuzzy coordinates of groundwater aquifers in Table 4 based on Algorithm 1, we have

$$\begin{aligned}
\hat{\theta} &= \frac{1}{3} [corr(x1, y1) + corr(x2, y2) + corr(x3, y3)] \\
&= \frac{1}{3} [corr((54.63, \dots, 55.50), (28.06, \dots, 28.60)) + corr((56.71, \dots, 57.10), (30.14, \dots, 30.33)) \\
&+ corr((58.78, \dots, 58.69), (32.21, \dots, 32.05))] \\
&= \frac{0.34 + 0.25 - 0.16}{3} = 0.14,
\end{aligned}$$

the value $\hat{\theta} = 0.14$ leads to the selection of the *FGM* copula to estimate the bivariate distribution using Algorithm 1. Therefore, by replacing Eqs. (13-14) in the *FGM* copula, the *cdf* estimation for (X, Y) is calculated by

$$\hat{F}_{\tilde{x}, \tilde{y}}(x, y) = \hat{F}_{\tilde{x}}(x) \hat{F}_{\tilde{y}}(y) \left(1 + 3\hat{\theta}(1 - \hat{F}_{\tilde{x}}(x))(1 - \hat{F}_{\tilde{y}}(y)) \right),$$

and also the *pdf* of (X, Y) is estimated by

$$\hat{f}_{\tilde{x}, \tilde{y}}(x, y) = \left(1 + 3\hat{\theta}(1 - 2\hat{F}_{\tilde{x}}(x))(1 - 2\hat{F}_{\tilde{y}}(y)) \right) \hat{f}_{\tilde{x}}(x) \hat{f}_{\tilde{y}}(y),$$

in which $\hat{f}_{\tilde{x}}(x) = \frac{1}{15} \sum_{i=1}^n \tilde{x}_i^*(x)$ is the estimation for univariate *PDF* of X based on fuzzy observations $\tilde{\mathbf{x}}$. The estimation of the bivariate *PDF* and *CDF* of geographical coordinates for the place of groundwater aquifers is plotted in Figure 3. Moreover, the contour plot of the bivariate *CDF* for (X, Y) with its marginal distributions is depicted in Figure 4.

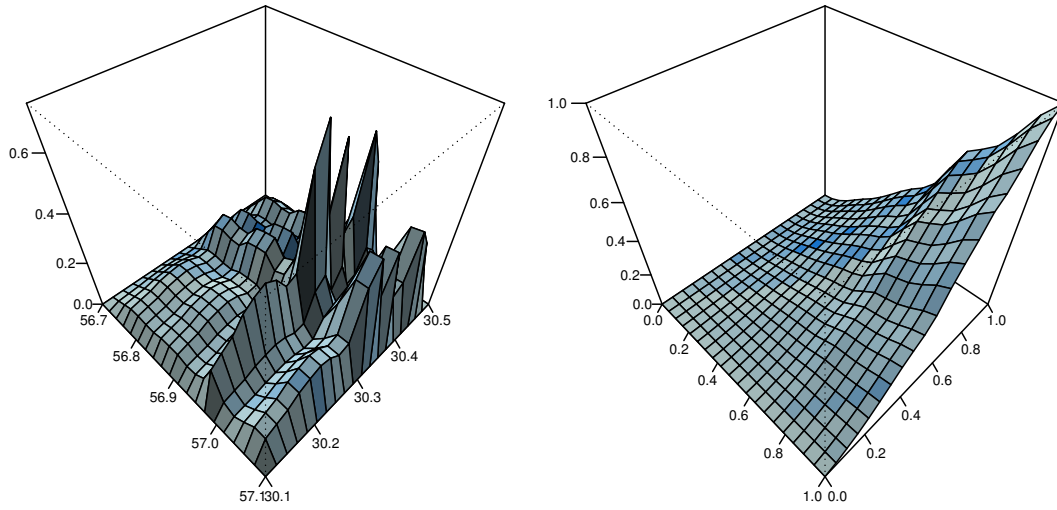


Figure 3: The estimation of bivariate *pdf* and *cdf* for groundwater aquifers based on *FGM* copula with $k = 6$.

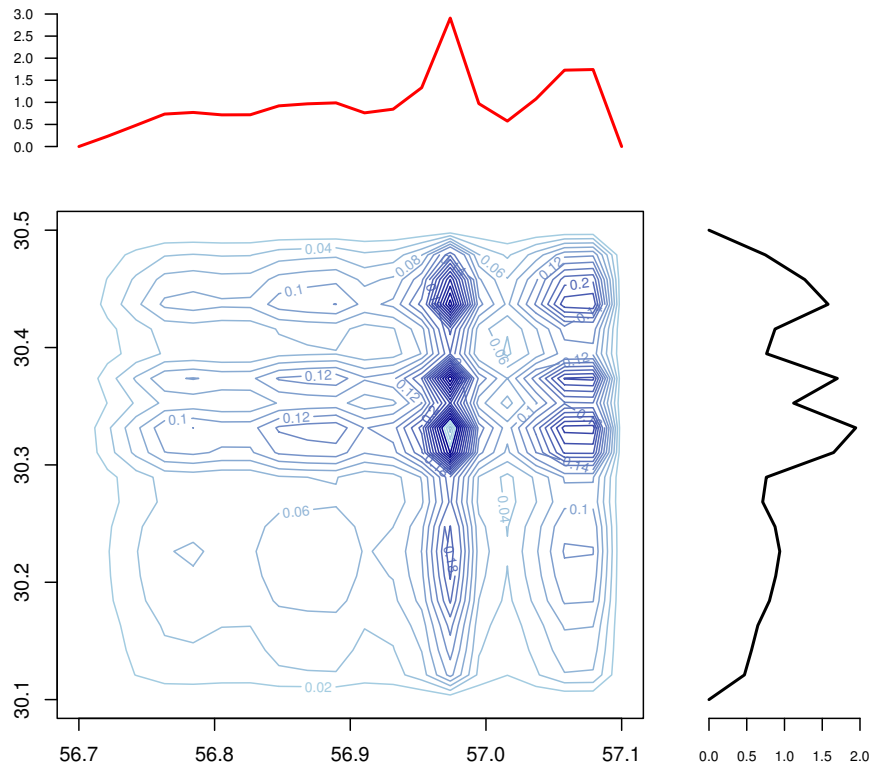


Figure 4: Contour plot for bivariate *pdf* estimation of groundwater aquifers (blue contour lines) based on *FGM* copula with $k = 6$, along with its marginal distributions: the univariate *PDF* of X (red curve) at the top and the univariate *PDF* of Y (black curve) on the right.

The results of this modeling can be utilized in water resource planning and management. By having the bivariate

distribution of groundwater aquifers, it becomes possible to identify high-yield areas and determine optimal locations for drilling new wells. Moreover, this approach provides a better understanding of the geographical spread of groundwater aquifers, enabling more informed and effective decision-making regarding their utilization.

6 Conclusions and future works

In this paper, a new approach for modeling bivariate distributions by fuzzy observations is presented. This method combines copula functions with the univariate density estimation method. Many traditional statistical methods face challenges when dealing with the uncertainty and vagueness inherent in fuzzy data. The proposed approach focuses on calculating the correlation coefficient between variables with fuzzy observations, where the appropriate copula is selected based on the strength of the correlation to link the marginal distributions. This process is carried out using Sklar's theorem which allows for the selection and combination of marginal distributions with copula functions. Moreover, the proposed modeling approach is applied to estimate the distribution of groundwater aquifers where their geographical coordinates are fuzzy rather than crisp.

Open problem A fundamental challenge for future research is establishing the consistency of estimators for the probability distribution function constructed from fuzzy data. Specifically, as the sample size increases, does the cumulative distribution function derived from fuzzy data converge to the true underlying cumulative distribution function? In other words, if $\hat{F}_{\tilde{x}}$ is an estimator based on a sample of n fuzzy observations, under what conditions do we have

$$\hat{F}_{\tilde{x}}(x) \xrightarrow{P} F(x),$$

(convergence in probability) or even

$$\sup_x |\hat{F}_{\tilde{x}}(x) - F(x)| \xrightarrow{\text{a.s.}} 0,$$

(almost sure uniform convergence)

Similarly, in the bivariate model, if $\hat{F}_{\tilde{x},\tilde{y}}$ is an estimator based on a sample of n fuzzy bivariate observations, under what conditions do we have

$$\hat{F}_{\tilde{x},\tilde{y}}(x, y) \xrightarrow{P} F(x, y),$$

or even

$$\sup_{(x,y) \in \mathbb{R}^2} |\hat{F}_{\tilde{x},\tilde{y}}(x, y) - F(x, y)| \xrightarrow{\text{a.s.}} 0.$$

Further theoretical investigation on imprecise data simulation is needed to determine the necessary and sufficient conditions for such convergence in both univariate and bivariate cases.

Acknowledgement

The authors would like to thank the respected anonymous referees whose comments improved the paper.

References

- [1] M. Arefi, R. Viertl, S. M. Taheri, *Fuzzy density estimation*, *Metrika*, **75** (2012), 5-22. <https://doi.org/10.1007/s00184-010-0311-y>
- [2] W. Breymann, A. Dias, P. Embrechts, *Dependence structures for multivariate high-frequency data in finance*, *Quantitative Finance*, **3**(1) (2003), 1-14. <https://doi.org/10.1080/713666155>
- [3] U. Cherubini, E. Luciano, W. Vecchiato, *Copula methods in finance*, John Wiley and Sons, **2** (2004), 949-956. <https://doi.org/10.1002/9781118673331>
- [4] G. Hesamian, M. G. Akbari, *Nonparametric kernel estimation based on fuzzy random variables*, *IEEE Transactions on Fuzzy Systems*, **25**(1) (2016), 84-99. <https://doi.org/10.1109/TFUZZ.2016.2551283>
- [5] J. S. Huang, S. Kotz, *Modifications of the Farlie-Gumbel-Morgenstern distributions. A tough hill to climb*, *Metrika*, **49**(2) (1999), 135-145. <https://doi.org/10.1007/s001840050030>

- [6] H. Joe, *Multivariate models and multivariate dependence concepts*, Chapman and Hall/CRC Press, 1997. <https://doi.org/10.1201/9780367803896>
- [7] P. Khalilpoor, A. Parchami, R. Pourmousa, *Density estimation based on fuzzy data by inspiration of kernel estimation method*, Iranian Journal of Mathematical Sciences and Informatics, Accepted, (2025). <https://dx.doi.org/10.2139/ssrn.4801965>
- [8] B. Khorramdel, C. Y. Chung, N. Safari, G. C. D. Price, *A fuzzy adaptive probabilistic wind power prediction framework using diffusion kernel density estimators*, IEEE Transactions on Power Systems, **33**(6) (2018), 7109-7121. <https://dx.doi.org/10.1109/tpwrs.2018.2848207>
- [9] D. Kim, J. M. Kim, S. M. Liao, Y. S. Jung *Mixture of D-vine copulas for modeling dependence*, Computational Statistics and Data Analysis, **64** (2013), 1-19. <https://dx.doi.org/10.1016/j.cstda.2013.02.018>
- [10] P. Kumar, *Probability distributions and estimation of Ali-Mikhail-Haq copula*, Applied Mathematical Sciences, **4**(14) (2010), 657-666.
- [11] C. D. Lai, N. Balakrishnan, *Continuous bivariate distributions*, Springer-Verlag New York, 2009. https://doi.org/10.1007/b101765_15
- [12] P. A. Lewis, *Distribution of the Anderson-Darling statistic*, The Annals of Mathematical Statistics, **32**(4) (1961), 1118-1124. <https://doi.org/10.1214/aoms/1177704850>
- [13] S. Lukasik, P. A. Kowalski, M. Charytanowicz, P. Kulczycki, *Fuzzy models synthesis with kernel-density-based clustering algorithm*, In 2008 Fifth International Conference on Fuzzy Systems and Knowledge Discovery, **3** (2008), 449-453. <https://doi.org/10.1109/FSKD.2008.139>
- [14] C. Meyer, *The bivariate normal copula*, Communications in Statistics-Theory and Methods, **42**(13) (2013), 2402-2422. <https://doi.org/10.1080/03610926.2011.611316>
- [15] M. R. Mortuza, E. Moges, Y. Demissie, H. Y. Li, *Historical and future drought in Bangladesh using copula-based bivariate regional frequency analysis*, Theoretical and Applied Climatology, **135** (2019), 855-871. <https://doi.org/10.1007/s00704-018-2407-7>
- [16] R. B. Nelsen, *An introduction to copulas*, Springer, 2006. <https://doi.org/10.1007/0-387-28678-0>
- [17] B. Renard, M. Lang, *Use of a Gaussian copula for multivariate extreme value analysis: Some case studies in hydrology*, Advances in Water Resources, **30**(4) (2007), 897-912. <https://doi.org/10.1016/j.advwatres.2006.08.001>
- [18] A. Shemyakin, A. Kniazev, *Introduction to Bayesian estimation and copula models of dependence*, John Wiley and Sons, 2017. <https://doi.org/10.1002/9781118959046>
- [19] M. A. Stephens, *EDF statistics for goodness of fit and some comparisons*, Journal of the American Statistical Association, **69**(347) (1974), 730-737. <https://doi.org/10.1080/01621459.1974.10480196>
- [20] X. Sun, M. Khayatnezhad, *Fuzzy-probabilistic modeling the flood characteristics using bivariate frequency analysis and α -cut decomposition*, Water Supply, **21**(8) (2021), 4391-4403. <https://doi.org/10.2166/ws.2021.186>
- [21] A. Urrutia, J. Galindo, L. Jimenéz, M. Piattini, *Data modeling dealing with uncertainty in fuzzy logic*, In: Avison, D., Elliot, S., Krogstie, J., Pries-Heje, J. (eds) The Past and Future of Information Systems: 1976-2006 and Beyond. IFIP WCC TC8 2006. IFIP International Federation for Information Processing, Vol 214. Springer, Boston, 2006. https://doi.org/10.1007/978-0-387-34732-5_19
- [22] R. Viertl, *Statistical methods for fuzzy data*, John Wiley and Sons, 2011. <https://doi.org/10.1002/9780470974414>