

# **A Goodness-of-fit Test for Semi-parametric Copula Models of Right-Censored Bivariate Survival Times**

by

**Moyan Mei**

B.Sc. (Honors), Dalhousie University, 2014

Project Submitted in Partial Fulfillment of the  
Requirements for the Degree of  
Master of Science

in the  
Department of Statistics and Actuarial Science  
Faculty of Science

© **Moyan Mei 2016**  
**SIMON FRASER UNIVERSITY**  
**Summer 2016**

All rights reserved.

However, in accordance with the *Copyright Act of Canada*, this work may be reproduced without authorization under the conditions for “Fair Dealing.” Therefore, limited reproduction of this work for the purposes of private study, research, criticism, review and news reporting is likely to be in accordance with the law, particularly if cited appropriately.

# Approval

**Name:** Moyan Mei  
**Degree:** Master of Science (Statistics)  
**Title:** *A Goodness-of-fit Test for Semi-parametric Copula Models of Right-Censored Bivariate Survival Times*  
**Examining Committee:** **Chair:** Tim Swartz  
Professor

**Michelle(Qian) Zhou**  
Senior Supervisor  
Assistant Professor

---

**Liangliang Wang**  
Supervisor  
Assistant Professor

---

**Rachel Altman**  
Internal Examiner  
Associate Professor

---

**Date Defended:** 9 June 2016

---

# Abstract

In multivariate survival analyses, understanding and quantifying the association between survival times is of importance. Copulas, such as Archimedean copulas and Gaussian copulas, provide a flexible approach of modeling and estimating the dependence structure among survival times separately from the marginal distributions (Sklar, 1959). However, misspecification in the parametric form of the copula function will directly lead to incorrect estimation of the joint distribution of the bivariate survival times and other model-based quantities.

The objectives of this project are two-folded. First, I reviewed the basic definitions and properties of commonly used survival copula models. In this project, I focused on semi-parametric copula models where the marginal distributions are unspecified but the copula function belongs to a parametric copula family. Various estimation procedures of the dependence parameter associated with the copula function were also reviewed. Secondly, I extended the pseudo in-and-out-of-sample (PIOS) likelihood ratio test proposed in Zhang et al. (2016) to testing the semi-parametric copula models for right-censored bivariate survival times. The PIOS test is constructed by comparing two forms of pseudo likelihoods, one is the "in-sample" pseudo likelihood, which is the full pseudo likelihood, and the other is the "out-of-sample" pseudo likelihood, which is a cross-validated pseudo likelihood by the means of jackknife. The finite sample performance of the PIOS test was investigated via a simulation study. In addition, two real data examples were analyzed for illustrative purposes.

**Keywords:** Goodness-of-fit; Right censoring data; Archimedean Copula; Gaussian Copula; Survival Analysis

# Contents

<b>Approval</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>Table of Contents</b>	<b>iv</b>
<b>Contents</b>	<b>iv</b>
<b>List of Tables</b>	<b>vi</b>
<b>List of Tables</b>	<b>vi</b>
<b>List of Figures</b>	<b>vii</b>
<b>List of Figures</b>	<b>vii</b>
<b>Acknowledgments</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Copula for Multivariate Survival Time . . . . .	2
1.1.1 Definition of Copula . . . . .	2
1.1.2 Survival copula . . . . .	2
1.2 Parametric Copula Families for Multivariate Survival Times . . . . .	2
1.2.1 Archimedean Copula . . . . .	2
1.2.2 Definition . . . . .	3
1.2.3 Gaussian family . . . . .	5
1.2.4 Dependence measurement . . . . .	5
1.3 Estimation of Dependence Structure . . . . .	10
1.4 Literature Review: Tests for Copula Models . . . . .	11
1.5 Organization and Objectives . . . . .	12

---

<b>2</b>	<b>In-and-Out-of-sample Pseudo Likelihood Ratio Test</b>	<b>13</b>
2.1	Definitions and notations . . . . .	13
2.2	Likelihood function . . . . .	14
2.3	Two-step Pseudo Maximum Likelihood Estimates of $\theta$ . . . . .	16
2.4	PIOS test . . . . .	18
<b>3</b>	<b>Numerical Illustration</b>	<b>20</b>
3.1	Simulation Studies . . . . .	20
3.1.1	Setups . . . . .	20
3.2	Application . . . . .	26
<b>4</b>	<b>Discussion</b>	<b>27</b>
	<b>Bibliography</b>	<b>29</b>
<b>A</b>	<b>Empirical standard error comparisons</b>	<b>33</b>

# List of Tables

Table 3.1	Type I error with sample size of 500 . . . . .	22
Table 3.2	Type I error with sample size of 1000 . . . . .	23
Table 3.3	Test power with sample size of 500 . . . . .	25
Table 3.4	Test power with sample size of 1000 . . . . .	25

# List of Figures

Figure 1.1	Plots of Gumbel Copula with $\theta$ at 1 and 8 . . . . .	8
Figure 1.2	Plots of Clayton Copula with $\theta = 8$ . . . . .	9
Figure 1.3	No tails dependence copulas . . . . .	9
Figure 3.1	Scatter plots of 500 pairs of simulated data from three different copulas with a common Kendall's $\tau = 0.8$ . . . . .	24

# Acknowledgments

Foremost I would like to express my deep gratitude to my advisor, Professor Michelle Zhou. Her inspiration, guidance, encouragement and insight helped me through my Master's study.

I am grateful to Professor Tim Swartz, Liangliang Wang, and Rachel Altman for serving on my examining committee and for their valuable comments and suggestions.

I also would like to give my heartfelt appreciation to the professors who I have taken courses with and the graduate students who I have studied with for their support and encouragement. This essay is dedicated to them.

Last but not least, a special acknowledgment goes to my mother and my girlfriend, who supported me throughout the whole writing process and kept providing endless encouragement.



# Chapter 1

## Introduction

Assessing dependency among multiple variables is a primary task in multivariate survival analysis. Copula models have appeared as a popular tool because they separate the joint distribution into two components: the marginal distributions of the individual variables, and the interdependency between the variables. Thus, the dependence structure can be handled separately from the marginal distributions, which provides great flexibility on the choices of both the marginal distributions and dependence structure. For a set of bivariate survival times  $(T_1, T_2)$  with respective marginal survival functions  $\mathbf{S}_1(t_1)$  and  $\mathbf{S}_2(t_2)$ , according to Sklar's theorem (Sklar, 1959), their joint survival distribution function  $\mathbf{S}(t_1, t_2)$  can be modeled in terms of univariate marginal distribution functions:

$$\mathbf{S}(t_1, t_2) = \mathbb{C}(\mathbf{S}_1(t_1), \mathbf{S}_2(t_2); \theta)$$

where  $\mathbb{C}(\cdot, \cdot; \theta)$  is a copula function with a parameter  $\theta$  controlling the dependence association between the survival times  $T_1$  and  $T_2$ . Usually, the marginal distributions of the survival times are unspecified. Shih and Louis (1995) proposed a semi-parametric estimator of the parameter of interest  $\theta$  via a two-step procedure. When a parametric copula model is used in applications, misspecification on its parametric structure may lead to inaccurate statistical estimation and inference. Thus, I am interested in the development of a goodness-of-fit test for misspecification in copula models of right-censored bivariate survival times.

---

## 1.1 Copula for Multivariate Survival Time

### 1.1.1 Definition of Copula

A  $d$ -dimensional copula,  $\mathbb{C}$ , is a  $d$ -dimensional distribution function with uniform marginals on  $[0, 1]$ . By Sklar's theorem (Sklar, 1959), for every cumulative distribution function (cdf)  $\mathbf{F}$ , of a  $d$ -dimensional continuous random vector  $\mathbf{X} = (X_1, \dots, X_d)$ , there exists a unique function  $\mathbb{C}$ , satisfying

$$\mathbf{F}(x_1, x_2, \dots, x_d) = \mathbb{C}(\mathbf{F}_1(x_1), \mathbf{F}_2(x_2), \dots, \mathbf{F}_d(x_d)),$$

where  $\mathbf{F}_j$  is the univariate marginal cdf of  $\mathbf{X}_j$ ,  $j = 1, 2, \dots, d$  and  $(x_1, x_2, \dots, x_d) \in \mathbb{R}^d$ . On the other hand, a copula function could be extracted from a multivariate distribution function. Given the joint cdf  $\mathbf{F}$  of a continuous random vector  $\mathbf{X}$  with univariate marginal cdfs  $\mathbf{F}_i$ 's, a copula function  $\mathbb{C}$  could be defined as:

$$\mathbb{C}(u_1, u_2, \dots, u_d) = \mathbf{F}(\mathbf{F}_1^{-1}(u_1), \mathbf{F}_2^{-1}(u_2), \dots, \mathbf{F}_d^{-1}(u_d)),$$

where  $(u_1, u_2, \dots, u_d) \in [0, 1]^d$  and  $\mathbf{F}_i^{-1}$ 's are the inverse of the marginal cdfs or marginal quantile functions.

### 1.1.2 Survival copula

Consider a  $d$ -dimensional multivariate survival times  $(T_1, T_2, \dots, T_d)$ . The multivariate survival function  $\mathbf{S}(t_1, t_2, \dots, t_d)$  can be expressed as a copula function of marginal survivals:

$$\mathbf{S}(t_1, t_2, \dots, t_d) = \mathbb{C}(\mathbf{S}_1(t_1), \mathbf{S}_2(t_2), \dots, \mathbf{S}_d(t_d)) \tag{1.1}$$

where  $\mathbf{S}_1(t_1), \mathbf{S}_2(t_2), \dots, \mathbf{S}_d(t_d)$  are marginal survival functions.

## 1.2 Parametric Copula Families for Multivariate Survival Times

### 1.2.1 Archimedean Copula

Archimedean copulas (Nelson, 2006) are the most commonly used copulas for modeling multivariate survival data. There are several reasons why the Archimedean copula family

---

is popular. First, Archimedean copula family consists of many parametric copula models, such as Clayton, Frank, Gumbel, independence, and other copulas, so this allows for a variety of dependence structures. Secondly, all commonly used Archimedean copulas have closed forms, unlike the copulas that are derived from multivariate distribution functions by Sklar's theorem. Thirdly, Archimedean copulas allow modeling dependence in high dimensions with only one parameter. This indicates that Archimedean copula is easy to handle and computationally straightforward. Last but not least, Shih and Louis (1995) states that Archimedean copula models correspond to proportional frailty models (Oakes, 1989) under certain circumstances. The main idea of proportional frailty models is to introduce dependence between survival times  $T_1$  and  $T_2$  by using an unobserved random variable  $G$ , the *frailty*, and here  $T_1$  and  $T_2$  are conditionally independent on  $G$ .

## 1.2.2 Definition

I begin with a general definition of Archimedean copulas, which can be found in Nelson (2006). Let  $\varphi_\theta$  be a continuous and strictly decreasing function from  $[0, 1]$  to  $[0, \infty]$  with a parameter  $\theta$ , where  $\varphi_\theta(1) = 0$ . The pseudo-inverse of  $\varphi_\theta$  is denoted by  $\varphi_\theta^{[-1]} : [0, \infty] \rightarrow [0, 1]$ , which is given by

$$\varphi_\theta^{[-1]}(t) = \begin{cases} \varphi_\theta^{-1}(t) & 0 \leq t \leq \varphi_\theta(0), \\ 0 & \varphi_\theta(0) \leq t \leq \infty. \end{cases}$$

Note that  $\varphi_\theta^{[-1]}$  is continuous and decreasing on  $[0, \infty]$ , and strictly decreasing on  $[0, \varphi_\theta(0)]$ . Furthermore,  $\varphi_\theta^{[-1]}(\varphi_\theta(u)) = u$  on  $[0, 1]$ , and

$$\varphi_\theta(\varphi_\theta^{[-1]}(t)) = \begin{cases} t & 0 \leq t \leq \varphi_\theta(0), \\ \varphi_\theta(0) & \varphi_\theta(0) \leq t \leq \infty. \end{cases}$$

Finally, if  $\varphi_\theta(0) = \infty$ , then  $\varphi_\theta^{[-1]} = \varphi_\theta^{-1}$ . A function  $\mathbb{C}$  from  $[0, 1]^2$  to  $[0, 1]$  given by

$$\mathbb{C}(u, v) = \varphi_\theta^{[-1]}(\varphi_\theta(u) + \varphi_\theta(v)), \quad (1.2)$$

is called an Archimedean copula. The function  $\varphi_\theta$  is called a generator of the copula, which is also the **inverse of a Laplace transform**. Furthermore, if  $\varphi_\theta(0) = \infty$ , then  $\varphi_\theta$  is a strict generator. In this case,  $\varphi_\theta^{[-1]} = \varphi_\theta^{-1}$ , and  $\mathbb{C}(u, v) = \varphi_\theta^{[-1]}(\varphi_\theta(u) + \varphi_\theta(v))$  is called a strict Archimedean copula. In the following, I introduce several commonly used Archimedean

---

copulas.

### **Gumbel family**

The Gumbel copula is given as

$$\mathbb{C}_\theta(u, v) = \varphi^{-1}(\varphi(u) + \varphi(v)) = \exp\left\{-\left[(-\log(u))^\theta + (-\log(v))^\theta\right]^{1/\theta}\right\},$$

where the generator function is  $\varphi_\theta(t) = (-\log(t))^\theta$  with  $\theta \geq 1$ . In addition,  $\varphi_\theta$  is a strict generator. It is easy to see that  $\varphi(t)$  is a continuous function of  $t$  and  $\varphi(1) = 0$ .

Note that, when  $\theta = 1$ , the Gumbel copula is an independent copula, i.e.,  $\mathbb{C}_\theta(u, v) = u \times v$ . When  $\theta \rightarrow \infty$ , the limit of Gumbel copula is a comonotonicity copula. The comonotonicity copula is the Fréchet–Hoeffding upper bound:  $M(u, v) = \min(u, v)$ , which corresponds to the perfectly positive dependence between two variables. Thus, the Gumbel copula interpolates between independence and perfectly positive dependence, and the parameter  $\theta$  reflects the strength of the dependence. For the Gumbel copula,  $\varphi_\theta$  is the Laplace form of a positive stable distribution. Hougaard (1986) stated that an important property of the use of stable distribution is that the proportionality of the conditional hazard given a frailty is inherited by the marginal hazard in univariate data.

### **Clayton family**

The Clayton copula is given as:

$$\mathbb{C}_\theta(u, v) = \max[(u^{-\theta} + v^{-\theta} - 1)^{-1/\theta}, 0],$$

where the generator function is  $\varphi_\theta(t) = (t^{-\theta} - 1)/\theta$  with  $\theta \in [-1, \infty)/\{0\}$ . For  $\theta > 0$ , the copula is strict and it is expressed as

$$\mathbb{C}_\theta(u, v) = (u^{-\theta} + v^{-\theta} - 1)^{-1/\theta}$$

As  $\theta \rightarrow 0$ , the Clayton copula approaches the independence copula, and as  $\theta \rightarrow \infty$ , it approaches the two-dimensional comonotonicity copula. Same as the Gumbel family, the Clayton copula also interpolates between independence and perfectly positive dependence. For the Clayton copula,  $\varphi_\theta$  is the Laplace form of gamma distribution. Clayton (1978)

---

shows that  $\lambda(t_2|T_1 = t_1)/\lambda(t_2|T_1 \geq t_1)$  equals to  $\theta$  if and only if the bivariate survival function belongs to the Clayton copula, where  $\lambda$  is the hazard function.

### Frank family

The Frank copula is given as:

$$\mathbb{C}_\theta(u, v) = -\frac{1}{\theta} \log \left( 1 + \frac{(e^{-\theta u} - 1)(e^{-\theta v} - 1)}{(e^{-\theta} - 1)} \right),$$

where the generator function is  $\varphi_\theta(t) = \log(e^{-\theta} - 1) - \log(e^{-\theta t} - 1)$  with  $\theta \neq 0$ . Frank copulas are strict Archimedean copulas. Furthermore, when  $\theta \rightarrow 0$ , the Frank copula approaches the independence copula; when  $\theta \rightarrow \infty$ , the Frank copula approaches the comonotonicity copula; when  $\theta \rightarrow -\infty$ , the Frank copula approaches the so-called the counter-comonotonicity copula, which corresponds to the perfectly negative dependence between two variables. Unlike the Gumbel and Clayton copulas, the Frank copula interpolates between perfectly negative dependence and perfectly positive dependence.

### 1.2.3 Gaussian family

In the literature, Archimedean copulas are most commonly used to model multivariate survival data, but they are not the only choices. The Gaussian copula was also considered, such as Li et al. (2008) and Othus and Li (2010). For  $\rho \in [-1, 1]$ , the Gaussian copula is defined as

$$\mathbb{C}(u, v) = \Phi_\rho(\Phi^{-1}(u), \Phi^{-1}(v)),$$

where

$$\Phi_\rho(x, y) = \int_{-\infty}^x \int_{-\infty}^y \frac{1}{2\pi\sqrt{1-\rho^2}} e^{\frac{2\rho st - s^2 - t^2}{2(1-\rho^2)}} ds dt,$$

and  $\Phi$  denotes the cdf of standard normal distribution. Same as the Frank copulas, the Gaussian copula also interpolates between perfectly negative dependence and perfectly positive dependence.

### 1.2.4 Dependence measurement

To quantify the dependence among two variables, several dependence measures have been proposed. Pearson's correlation is most frequently used in practice as a measure of dependence. Embrechts et al. (2000) shows that Pearson's correlation remains natural and

---

unproblematic in the class of elliptical distributions (e.g Gaussian and Student's  $t$ ) only. Thus, Pearson's correlation may not be appropriate to use if the distribution is not elliptical. In contrast to ordinary correlation measures, alternative measures of associations can be used, such as rank correlation and tail dependence.

### Kendall's tau

The Kendall's rank correlation, which is also called the Kendall's  $\tau$ , can be treated as a measure of concordance for bivariate random vectors. Let  $(x_1, y_1), (x_2, y_2)$  be two sets of observations of the bivariate random vector  $(X, Y)$ . Thus,  $(x_1, y_1), (x_2, y_2)$  are **concordant** if  $(x_2 - x_1)(y_2 - y_1) > 0$ , and  $(x_1, y_1), (x_2, y_2)$  are **discordant** if  $(x_2 - x_1)(y_2 - y_1) < 0$ .

It is easy to see that if  $Y$  tends to increase with  $X$ , then the probability of concordance is expected to be high relative to the probability of discordance. Otherwise, if  $Y$  tends to decrease with increasing  $X$ , the probability of concordance is expected to be low relative to the probability of discordance. This motivates Kendall's rank correlation, which is simply the probability difference between concordance and disconcordance for pairs of random vectors. Specifically, the Kendall's rank correlation is defined as

$$\rho_\tau(X, Y) = P[(x_2 - x_1)(y_2 - y_1) > 0] - P[(x_2 - x_1)(y_2 - y_1) < 0].$$

It can be expressed in a more compact way in the following

$$\rho_\tau(X, Y) = \mathbb{E}[\text{sign}((x_2 - x_1)(y_2 - y_1))],$$

where  $\text{sign}(x) = 1$  if  $x > 0$ ,  $= 0$  if  $x = 0$ , and  $= -1$  if  $x < 0$ . Note that:

- If the agreement between the two rankings is perfect (i.e., the two rankings are the same), then the coefficient is equal to 1.
- If the disagreement between the two rankings is perfect (i.e., one ranking is the reverse of the other), then the coefficient is equal to -1.
- If  $X$  and  $Y$  are independent, then  $\rho_\tau = 0$ .

Kendall's  $\tau$  can also be expressed in the form of the copula function  $\mathbb{C}_\theta(u, v)$  in the following:

$$\rho_\tau = -1 + 4 \int_0^1 \int_0^1 \mathbb{C}_\theta(u, v) c_\theta(u, v) du dv,$$

---

where  $c_\theta(u, v) = \partial \mathbb{C}_\theta(u, v) / \partial u \partial v$ . From this expression, we can see that the Kendall's  $\tau$  is just a function of  $\theta$ , and it only depends on the bivariate copula function, but not the marginal distributions.

For some of the copula families introduced above, Kendall's  $\tau$  is given as:

- Gumbel:  $\rho_\tau = 1 - 1/\theta$ ;
- Clayton:  $\rho_\tau = \theta/(\theta + 2)$ ;
- Gaussian:  $\rho_\tau = 2 \arcsin(\theta)/\pi$ .

### Tail Dependence

In some studies, the main interest focuses on the dependence between extreme values of the variables. Coefficients of tail dependence provide measures of the strength of dependence in the tails of a bivariate distribution.

Let  $X$  and  $Y$  be random variables with their respective cdfs  $\mathbf{F}_1$  and  $\mathbf{F}_2$ . The coefficient of upper tail dependence of  $X$  and  $Y$  is

$$\lambda_u(X, Y) = \lim_{q \rightarrow 1, q < 1} \Pr[Y > \mathbf{F}_2^-(q) | X > \mathbf{F}_1^-(q)],$$

provided the limit  $\lambda_u(X, Y) \in [0, 1]$  exists.  $\mathbf{F}_j^-$ s are the generalized inverse of marginal cdf  $\mathbf{F}_j$ s, where  $j = 1, 2$ . Mathematically, it is defined as  $\mathbf{F}_i^- = \inf\{x \in \mathbb{R} : \mathbf{F}(x) \geq y\}, y \in \mathbb{R}$ . If  $\lambda_u(X, Y) \in (0, 1]$ , then  $X$  and  $Y$  are said to show **upper tail dependence**. On the other hand, if  $\lambda_u(X, Y) = 0$ , they are **asymptotically independent** in the upper tail.

Analogously, the coefficient of lower tail dependence is defined as

$$\lambda_l(X, Y) = \lim_{q \rightarrow 1, q < 1} \Pr[Y \leq \mathbf{F}_2^-(q) | X \leq \mathbf{F}_1^-(q)]$$

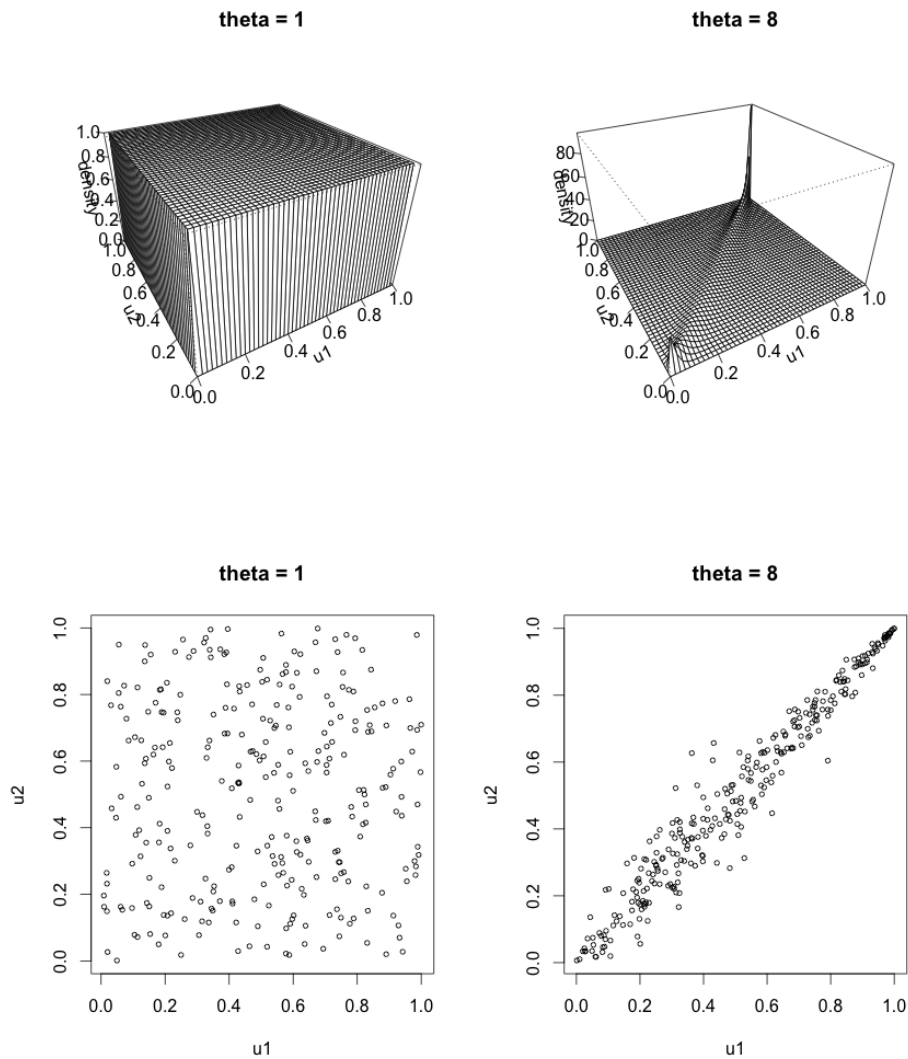
The coefficients of tail dependence can be expressed in the form of the copula functions

$$\lambda_u(X, Y) = \lim_{q \rightarrow 1, q < 1} \frac{\mathbb{C}(q, q)}{q} \quad \text{and} \quad \lambda_l(X, Y) = \lim_{q \rightarrow 1, q < 1} \frac{1 - 2q + \mathbb{C}(q, q)}{1 - q}.$$

Calculation of these coefficients is straightforward if the copula has an explicit form. For the Gumbel copula, the coefficient of upper tail dependence is  $\lambda_u(X, Y) = 2 - 2^{1/\theta}$ . Hence,

provided that  $\theta > 1$ , the Gumbel copula has upper tail dependence. The strength of this tail dependence approaches to 1 as  $\theta \rightarrow +\infty$ , while the strength of this tail dependence approaches to 0 as  $\theta \rightarrow 1$ . Figure 1.1 is an illustration of Gumbel copula with  $\theta$  equals to 1 and 8. The top two plots are the density perspective plots, and the bottom two are the scatter plots of simulated observations from the Gumbel copulas. Top left plot in Figure 1.1 demonstrates that there is no tail dependence at all, and the data seems random overall, which is supported by the bottom left scatter plot. In contrast, when  $\theta$  increases to 8, it shows an obvious upper tail in the top right plot. The scatter plot in the bottom right also suggests that there is a strong association between variables in the upper tail.

Figure 1.1: Plots of Gumbel Copula with  $\theta$  at 1 and 8





Unlike the Gumbel copula, the Clayton copula has lower tail dependence. The corresponding coefficient of lower tail dependence is  $\lambda_l(X, Y) = 2^{-1/\theta}$ , where  $\theta > 0$ . Figure 1.2 has shown that the two variables appear to behave closely in the left corner.

Figure 1.2: Plots of Clayton Copula with  $\theta = 8$

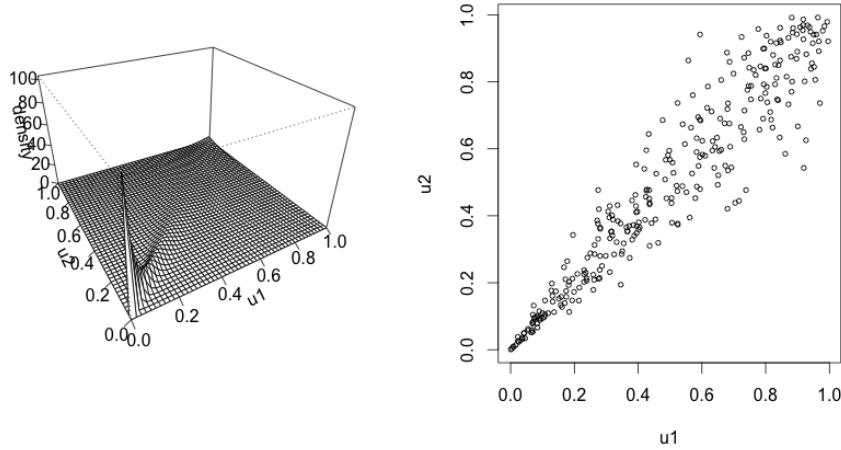
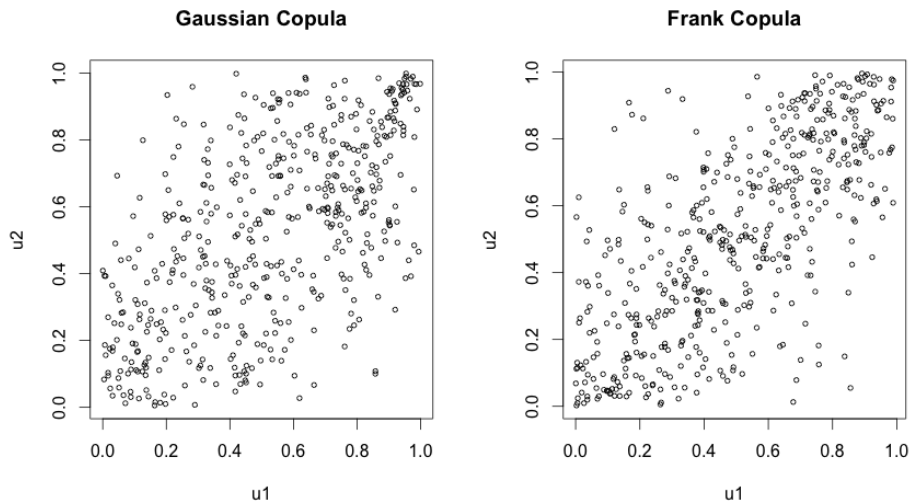


Figure 1.3: No tails dependence copulas



However, there exists some copulas that do not exhibit tail dependency. The Frank copula does not show dependency in both tails for any values of  $\theta$ . Moreover, the bivariate Gaussian copula is also symmetric. Please see Figure 1.3.

---

### 1.3 Estimation of Dependence Structure

In the literature, several procedures have been proposed to estimate the dependence parameter of the copula function. When both the marginal distributions and the copula functions are parametrically specified, the maximum likelihood estimates of the parameters associated with the marginal distributions and the copula can be obtained simultaneously. Alternatively, a two-stage estimation procedure (Shih and Louis, 1995) can be used. In this procedure, the parameters associated with the marginal distributions and the parameters associated with the copula function are estimated separately in two steps.

In practice, however, the true marginal distributions are rarely known. A recent empirical study in Kim et al. (2007) indicates that the fully parametric method is not robust against misspecification of the marginals. Thus, many authors advocate the use of a two-stage semi-parametric estimation procedure studied in Genest et al. (1995) and Shih and Louis (1995) to estimate the dependence parameter  $\theta$ . In the first stage, the marginal survival functions are estimated by a non-parametric method, and in the second stage, the dependence parameter  $\theta$  is estimated by maximizing the pseudo-likelihood function with the estimated marginal distributions. Without censoring, Genest et al. (1995) estimated the marginal survival functions by the empirical distribution functions; with censoring, Shih and Louis (1995) used the Kaplan-Meier estimators. Shih and Louis (1995) also investigated the asymptotic properties of the two-stage pseudo maximum likelihood estimators (PMLE). They showed that  $\sqrt{n}(\widehat{\theta}-\theta^*)$  converges in distribution to a normal random variable with mean zero under some regularity conditions, where  $\widehat{\theta}$  is the PMLE of the parameter  $\theta$ , and  $\theta^*$  is the limiting value of  $\widehat{\theta}$ .

Non-parametric estimation approaches of the copula functions have also been considered. Gijbels and Mielniczuk (1990) and Chen and Huang (2007) proposed kernel estimators of the copulas given uncensored data. Chen and Huang (2007) formed their nonparametric copula estimator in two stages. First, kernel estimators of the marginal distribution functions are obtained. In the second stage, a kernel copula estimator is obtained based on local linear kernels and a simple mathematical correction that removes the boundary bias. Given censored data, Gribkova and Lopez (2015) defined one discrete and two smooth estimators of the copula. The first smooth estimator generalizes the estimator of Fermanian et al. (2004) to the censored survival times, while the second estimator followed the method proposed by Omelka et al. (2009) with correction. Moreover, the convergence rates of both estimators were derived. Other work of non-parametric estimation of copula functions can

---

be found in Kalbfleisch and Prentice (2002), Tsai et al. (1997), Lin and Ying (1993), Prentice et al. (1997), and Prentice et al. (2004).

## 1.4 Literature Review: Tests for Copula Models

Various goodness-of-fit tests for copulas have been proposed throughout the literature. Malevergne and Sornette (2003) developed a chi-square test for testing the Gaussian copulas. Fermanian (2005) introduced two distribution-free goodness-of-fit test statistics for copulas. The first test used the chi-square test, and the observed and expected frequencies considered based on a kernel estimate of the copula density and a parametric estimate of the copula density respectively. The second test was based on the minimum distance between the smoothed copula density and the estimated parametric density. Scaillet (2007) also used a kernel based goodness-of-fit test for copulas with fixed smoothing parameters. Additionally, there are other types of specification tests. Prokhorov and Schmidt (2009) considered a conditional moment test of the validity of the copula. Mesfioui et al. (2009) proposed a test based on Spearman dependence function. Genest et al. (2011) came up with a Cramér–von Mises type statistic, which compares the distance between an estimate of the parametric Pickands dependence function and nonparametric estimators studied by Genest and Segars (2009). More recently, Huang and Prokhorov (2014) provided a rank-based goodness-of-fit test for copulas based on the information test by White (1982).

However, these tests were designed for fully observed data, and they cannot be applied to the censored data. Several authors have proposed several goodness-of-fit tests for Archimedean copulas of censored bivariate survival times. Shih (1998) proposed a goodness-of-fit test procedure for the bivariate Clayton model. The test compares unweighted and weighted concordance estimators of the association parameter  $\theta$ , and if the assumed Clayton model is true, these estimators converge to the true value of  $\theta$  and their difference should be close to zero. Wang and Wells (2000) proposed a model selection procedure for right censored bivariate data based on the  $L_2$  norm of the Kendall's distribution, which is basically the measurement of the distance between the empirical and model-based estimates of Kendall's distribution. Genest et al. (2006) extended the idea of Wang and Wells (2000) to propose a goodness-of-fit test based on the probability integral transformation, and it also offered a way to compute the asymptotic  $p$  value for various goodness-of-fit tests based on a non-truncated Kendall's process. Andersen et al. (2005) proposed three test statistics to check whether an assumed one-parameter shared frailty model fits bivariate right censored

---

data without covariates. The idea behind Andersen et al. (2005) is to measure the difference between the semi-parametric estimate and the non-parametric estimate of a proposed copula model via chi-squared type statistic, Kolmogorov-like statistic and weighted difference based statistic. Later, Emura et al. (2010) extended the idea of comparison between two point estimators derived under the same class of estimation equations with different weight functions to general Archimedean models.

In this paper, I considered an alternative test procedure, the pseudo in-and-out-of sample (PIOS) likelihood ratio test in Zhang et al. (2016), for semi-parametric copulas on right-censored bivariate survival times. The idea of PIOS statistic is rooted in the in-and-out-of-sample (IOS) likelihood ratio test by Presnell and Boos (2004) for cross-sectional univariate data. The IOS test provided a measure on how sensitive the likelihood is to the varying data by the means of jackknife. The PIOS test is asymptotical equivalent to the information ratio (IR) test originally proposed by Zhou et al. (2012). Later, both the PIOS test and IR test have been extended to test the model mis-specifications for univariate and multivariate time series data by Zhang et al. (2016). However, those tests are designed for fully observed data. In this project, I applied the PIOS test in testing the semi-parametric copula models for right censored bivariate survival times. Compared with the fully observed data, the likelihood function for censored data is more complicated. This results in more challenges in estimating the marginal distributions including the marginal survival functions and density functions.

## **1.5 Organization and Objectives**

This paper is organized as follows. In Chapter 2, I will introduce the statistical procedure of testing semi-parametric copula models for right-censored bivariate survival times. In Chapter 3, I will investigate the performance of the PIOS test through simulation study and two real data examples. In the last chapter, I will conclude my thesis with discussions and future work.

# Chapter 2

## In-and-Out-of-sample Pseudo Likelihood Ratio Test

### 2.1 Definitions and notations

First, I introduce some notations used in the following sections. Suppose  $T_{i1}$  and  $T_{i2}$  are the  $i$ -th bivariate event times, and  $C_i$  is the  $i$ -th censoring time, where  $i = 1, 2, \dots, n$ . Let  $X_{i1}$  and  $X_{i2}$  denote the observed times, where  $X_{ij} = \min(T_{ij}, C_i)$ , for  $j = 1, 2$ . Let  $\delta_{ij} = I(T_{ij} \leq C_i)$ ,  $j = 1, 2$  denote the censoring indicator variables. In summary, the observed data can be expressed as  $\mathfrak{D} = \{(X_{i1}, X_{i2}, \delta_{i1}, \delta_{i2}), i = 1, 2, \dots, n\}$ . Here, I assume independent censoring, i.e.  $C_i$  is independent of the bivariate event times  $T_{i1}$  and  $T_{i2}$ .

Let  $\mathbf{S}(t_1, t_2) = \Pr(T_{i1} > t_1, T_{i2} > t_2)$  be the joint survival function of  $(T_{i1}, T_{i2})$ , and let  $\mathbf{S}_1(t_1) = \Pr(T_{i1} > t_1)$  and  $\mathbf{S}_2(t_2) = \Pr(T_{i2} > t_2)$  be the marginal survival functions of  $T_{i1}$  and  $T_{i2}$  respectively. We assume that there exists a copula function  $\mathbb{C}_0$  such that

$$\mathbf{S}(t_1, t_2) = \Pr(T_{i1} > t_1, T_{i2} > t_2) = \mathbb{C}_0(\mathbf{S}_1(t_1), \mathbf{S}_2(t_2); \theta),$$

where  $\mathbb{C}_0$  is the true copula. However, it is rarely known. Thus, we model the joint survival function  $\mathbf{S}(t_1, t_2)$  by a copula family in the following form

$$\mathbf{S}(t_1, t_2) = \mathbb{C}(\mathbf{S}_1(t_1), \mathbf{S}_2(t_2); \theta), \tag{2.1}$$

where  $\mathbb{C}$  is a copula function, and  $\theta$  is the dependence parameter. As mentioned earlier, the Archimedean copulas and the Gaussian copula are usually considered to model multivariate survival data in the literature, and both of these two copula families are specified by one

---

dependence parameter. Thus, in this thesis, for the purpose of simplicity, I only consider copula families with one dependence parameter.

The goal of this paper is to propose a test for

$$H_0 : \mathbb{C}_0 \in C = \{\mathbb{C}(\cdot; \theta), \theta \in \Theta\} \quad \text{vs.} \quad H_A : \mathbb{C}_0 \notin C = \{\mathbb{C}(\cdot; \theta), \theta \in \Theta\},$$

where  $\Theta \subset \mathbb{R}$  is the parameter space.

In addition, some related notations are defined as follows. Let

- $f_1(t_1) = -\mathbf{S}'_1(t_1)$  and  $f_2(t_2) = -\mathbf{S}'_2(t_2)$  be the marginal density functions of  $T_{i1}$  and  $T_{i2}$  respectively;
- $U_{i1} = \mathbf{S}_1(X_{i1})$  and  $U_{i2} = \mathbf{S}_2(X_{i2})$ ;
- $\mathbb{C}_\theta(u_1, u_2) = \mathbb{C}(u_1, u_2; \theta)$
- $\mathbb{c}_\theta(u_1, u_2) = \partial \mathbb{C}_\theta(u_1, u_2) / \partial u_1 \partial u_2$
- $\mathbb{c}_\theta^{(j)}(u_1, u_2) = \partial \mathbb{C}_\theta(u_1, u_2) / \partial u_j, j = 1, 2.$

## 2.2 Likelihood function

Given the data  $\mathfrak{D}$ , there are generally four different types of data:

- No censoring for neither  $T_{i1}$  nor  $T_{i2}$ , so  $\delta_{i1} = \delta_{i2} = 1$ , and the corresponding likelihood component is  $\partial \mathbf{S}(X_{i1}, X_{i2}) / \partial X_{i1} \partial X_{i2}$
- Only  $T_{i1}$  is censored, but  $T_{i2}$  is observed, so  $\delta_{i1} = 0$  and  $\delta_{i2} = 1$ , and the corresponding likelihood component is  $-\partial \mathbf{S}(X_{i1}, X_{i2}) / \partial X_{i2}$
- Only  $T_{i2}$  is censored, but  $T_{i1}$  is observed, so  $\delta_{i1} = 1$  and  $\delta_{i2} = 0$ , and the corresponding likelihood component is  $-\partial \mathbf{S}(X_{i1}, X_{i2}) / \partial X_{i1}$
- Both  $T_{i1}$  and  $T_{i2}$  are censored, so  $\delta_{i1} = \delta_{i2} = 0$ , and the corresponding likelihood component is  $\mathbf{S}(X_{i1}, X_{i2})$

Thus, the full log-likelihood is written as  $\ell_{\mathfrak{D}}(\theta) = \sum_{i=1}^n \ell(\theta; \mathfrak{D}_i)$ , where  $\mathfrak{D}_i = (X_{i1}, X_{i2}, \delta_{i1}, \delta_{i2})$ , and the  $i$ -th component of the log-likelihood is

$$\begin{aligned} \ell(\theta; \mathfrak{D}_i) &= \delta_{i1}\delta_{i2} \log \left\{ \frac{\partial \mathbf{S}(X_{i1}, X_{i2})}{\partial X_{i1} \partial X_{i2}} \right\} + (1 - \delta_{i1})\delta_{i2} \log \left\{ -\frac{\partial \mathbf{S}(X_{i1}, X_{i2})}{\partial X_{i2}} \right\} \\ &\quad + (1 - \delta_{i2})\delta_{i1} \log \left\{ -\frac{\partial \mathbf{S}(X_{i1}, X_{i2})}{\partial X_{i1}} \right\} + (1 - \delta_{i1})(1 - \delta_{i2}) \log \{ \mathbf{S}(X_{i1}, X_{i2}) \}. \end{aligned} \quad (2.2)$$

Based on the survival copula in (2.1), the  $i$ -th log-likelihood component can be written as

$$\begin{aligned} \ell(\theta; \mathfrak{D}_i) &= \delta_{i1}\delta_{i2} \log \{ \mathbb{C}_{\theta}(U_{i1}, U_{i2}) f_1(X_{i1}) f_2(X_{i2}) \} + (1 - \delta_{i1})\delta_{i2} \log \left\{ -\frac{\partial \mathbb{C}_{\theta}(U_{i1}, U_{i2})}{\partial U_{i2}} \frac{\partial U_{i2}}{\partial X_{i2}} \right\} \\ &\quad + (1 - \delta_{i2})\delta_{i1} \log \left\{ -\frac{\partial \mathbb{C}_{\theta}(U_{i1}, U_{i2})}{\partial U_{i1}} \frac{\partial U_{i1}}{\partial X_{i1}} \right\} + (1 - \delta_{i1})(1 - \delta_{i2}) \log \{ \mathbb{C}_{\theta}(U_{i1}, U_{i2}) \} \end{aligned} \quad (2.3)$$

Note that all the Archimedean copulas and the Gaussian copula can be used in equation (2.3). For example, for the Gumbel family,

$$\mathbb{C}_{\theta}(u, v) = \varphi^{-1}(\varphi(u) + \varphi(v)) = \exp \left\{ -[(-\log(u))^{\theta} + (-\log(v))^{\theta}]^{\frac{1}{\theta}} \right\},$$

in the  $i$ -th log-likelihood components, it is given as

$$\begin{aligned} \frac{\partial \mathbf{S}(X_{i1}, X_{i2})}{\partial X_{i1}} &= \frac{\partial \mathbb{C}_{\theta}(U_{i1}, U_{i2})}{\partial U_{i1}} \times \frac{\partial U_{i1}}{\partial X_{i1}} \\ &= \mathbb{C}_{\theta}(U_{i1}, U_{i2}) \frac{\partial}{\partial U_{i1}} \left\{ -[(-\log(U_{i1}))^{\theta} + (-\log(U_{i2}))^{\theta}]^{1/\theta} \right\} \frac{\partial U_{i1}}{\partial X_{i1}} \\ &= \mathbb{C}_{\theta}(U_{i1}, U_{i2}) \left\{ -1/\theta [(-\log(U_{i1}))^{\theta} + (-\log(U_{i2}))^{\theta}]^{(1-\theta)/\theta} \right\} \\ &\quad \times (\theta(-\log(U_{i1}))^{\theta-1})(-1/U_{i1})(-f_1(X_{i1})) \\ &= \mathbb{C}_{\theta}(U_{i1}, U_{i2}) [(-\log(U_{i1}))^{\theta} + (-\log(U_{i2}))^{\theta}]^{(1-\theta)/\theta} (-\log(U_{i1}))^{\theta-1} \frac{-f_1(X_{i1})}{U_{i1}} \end{aligned}$$

Similarly, the partial derivative with respect to  $X_{i2}$  can be obtained as follows:

$$\begin{aligned} \frac{\partial \mathbf{S}(U_{i1}, U_{i2})}{\partial X_{i2}} &= \frac{\partial \mathbb{C}_{\theta}(U_{i1}, U_{i2})}{\partial U_{i2}} \times \frac{\partial U_{i2}}{\partial X_{i2}} \\ &= \mathbb{C}_{\theta}(U_{i1}, U_{i2}) [(-\log(U_{i1}))^{\theta} + (-\log(U_{i2}))^{\theta}]^{(1-\theta)/\theta} (-\log(U_{i2}))^{\theta-1} \frac{-f_2(X_{i2})}{U_{i2}} \end{aligned}$$

At last, for the case where  $\delta_1 = \delta_2 = 1$ , I have

$$\begin{aligned} \frac{\partial^2 \mathbf{S}(U_{i1}, U_{i2})}{\partial X_{i1} \partial X_{i2}} &= c_\theta(U_{i1}, U_{i2}) f_1(X_{i1}) f_2(X_{i2}) \\ &= \mathbb{C}_\theta(U_{i1}, U_{i2}) (U_{i1} U_{i2})^{-1} [(-\log(U_{i1}))^\theta + (-\log(U_{i2}))^\theta]^{(2/\theta-2)} (\log(U_{i1}) \log(U_{i2}))^{\theta-1} \\ &\quad \times \left\{ 1 + (\theta - 1) [(-\log(U_{i1}))^\theta + (-\log(U_{i2}))^\theta]^{-1/\theta} \right\} f_1(X_{i1}) f_2(X_{i2}) \end{aligned}$$

### 2.3 Two-step Pseudo Maximum Likelihood Estimates of $\theta$

In this thesis, I consider situations where the marginal distributions of  $T_{i1}$  and  $T_{i2}$  are unspecified. A two-step estimation procedure was proposed in Shih and Louis (1995). In the first step, the marginal survival functions of  $T_{i1}$  and  $T_{i2}$  can be estimated non-parametrically via the Kaplan-Meier estimator (Kaplan and Meier, 1958) or Nelson-Aalen estimator (Nelson 1972 and Aalen 1978). Specifically, for the  $j$ -th event time,  $j = 1, 2$ , let  $x_{(1),j}, x_{(2),j}, \dots, x_{(k_j),j}$  be the distinct, ordered and uncensored event times. The Kaplan-Meier estimate of  $\mathbf{S}_j(t)$  is given by

$$\widehat{\mathbf{S}}_j(t) = \prod_{x_{(m),j} < t} \left[ 1 - \frac{dN(x_{(m),j})}{Y(x_{(m),j})} \right],$$

where  $Y(x_{(m),j}) = \sum_{i=1}^n I(X_{ij} > x_{(m),j})$  is total number of subjects at risk at time  $x_{(m),j}$ , and  $dN(x_{(m),j}) = \sum_{i=1}^n I(X_{ij} \leq x_{(m),j}) I(\delta_{ij} = 1)$  is total number of uncensored events prior to time  $x_{(m),j}$ ,  $m = 1, 2, \dots, k_j$ . In R, Kaplan-Meier estimators are obtained via the function `survfit` in package `survival`. The alternative way to estimate survival functions is to use **Nelson-Aalen estimators**. It is given as  $\widetilde{\mathbf{S}}_j(t) = \exp\{-\widetilde{\Lambda}_j(t)\}$ , where

$$\widetilde{\Lambda}_j(t) = \sum_{x_{(m),j} < t} \frac{dN(x_{(m),j})}{Y(x_{(m),j})}$$

is the estimate of the cumulative hazard function  $\Lambda_j(t)$ . The Kaplan-Meier estimator and the Nelson-Aalen estimator are asymptotically equivalent.



In the second step, the dependence parameter  $\theta$  of a copula function is estimated by maximizing the pseudo likelihood function

$$\begin{aligned} \widehat{\ell}(\theta; \mathfrak{D}) = & \sum_{i=1}^n \delta_{i1} \delta_{i2} \log \left\{ \mathbb{C}_{\theta}(\widehat{U}_{i1}, \widehat{U}_{i2}) \widehat{f}_1(X_{i1}) \widehat{f}_2(X_{i2}) \right\} + (1 - \delta_{i1}) \delta_{i2} \log \left\{ \frac{\partial \mathbb{C}_{\theta}(\widehat{U}_{i1}, \widehat{U}_{i2})}{\partial X_{i2}} \frac{\partial \widehat{U}_{i2}}{\partial X_{i2}} \right\} \\ & + (1 - \delta_{i2}) \delta_{i1} \log \left\{ \frac{\partial \mathbb{C}_{\theta}(\widehat{U}_{i1}, \widehat{U}_{i2})}{\partial X_{i1}} \frac{\partial \widehat{U}_{i1}}{\partial X_{i1}} \right\} + (1 - \delta_{i1})(1 - \delta_{i2}) \log \left\{ \mathbb{C}_{\theta}(\widehat{U}_{i1}, \widehat{U}_{i2}) \right\}, \end{aligned}$$

where  $\widehat{U}_{i1} = \widehat{\mathbf{S}}_1(X_{i1})$  and  $\widehat{U}_{i2} = \widehat{\mathbf{S}}_2(X_{i2})$  are the estimated marginal functions from the first step.

To construct the pseudo likelihood function, it is required to obtain the estimates of the density functions,  $\widehat{f}_1$  and  $\widehat{f}_2$ . Here, I suggest using kernel smoothed estimates, denoted by  $\widehat{\lambda}_j(t)$ , of the hazard functions  $\lambda_j(t) = \frac{f_j(t)}{\mathbf{S}_j(t)}$  (Proschan, 1963) to obtain the estimates of the density functions by  $\widehat{f}_j(t) = \widehat{\lambda}_j(t) \widehat{\mathbf{S}}_j(t)$ . More specifically,

$$\widehat{\lambda}_j(t) = \frac{1}{b(t)} \sum_{i=1}^n K\left(\frac{t - x_{(i),j}}{b(t)}\right) \frac{\delta_{(i),j}}{n - i + 1}.$$

where  $K(\cdot)$  is a kernel function,  $b(t)$  is the bandwidth of a kernel function that controls the smoothness of the estimated function,  $0 < x_{(1),j} < x_{(2),j} < \dots < x_{(n),j}$  are ordered times of  $\{X_{ij}, i = 1, \dots, n\}$  and  $\delta_{(i),j}$  is the corresponding indicator variable of  $x_{(i),j}$ . In R, the function `muhaz` in package `muhaz` can be used to obtain the estimate described above. Kernel smoothing is widely used in statistical applications, particularly for density functions and regression models. The performance of kernel smoothed estimates depends on the choice of bandwidth. A large bandwidth would lead to estimators with less variability at the cost of increased bias and inaccurate inference by the re-sampling. On the other hand, a small bandwidth would yield estimators with large variability. To find the optimal bandwidth for the kernel smoothed estimate of the hazard function, one can use locally optimal bandwidth selection method. There are two main advantages of locally optimal bandwidth method. First, an optimal local bandwidth can be consistently estimated by minimizing an estimate of the local mean squared error of the hazard rate estimate with respect to the bandwidth. Secondly, the local bandwidth estimates allow for larger bandwidths at points with larger variance and generally lead to increasing bandwidths near the left endpoint and towards the right endpoint. More details about the local bandwidth can be found in Wang et al. (1998). There are two popular methods to choose the optimal local bandwidth, which are

cross-validation and "plug-in" techniques respectively. In R, the optimal local bandwidth is obtained by "plug-in" technique through the function `muhaz` with argument `bw.method = "local"`. However, the `muhaz` function can only return the hazard estimates on a grid of ordered time events. Thus, we use `approx` function of R to obtain the hazard estimates on the observed event time in the data by interpolation.

Finally, a pseudo maximum likelihood estimates (PMLE) of  $\theta$  can be obtained by solving the equation

$$\frac{\partial \widehat{\ell}(\theta; \mathfrak{D})}{\partial \theta} = 0.$$

## 2.4 PIOS test

Essentially, the construction of the PIOS test consists of two parts: "in-sample" pseudo log likelihood and the "out-of-sample" pseudo log likelihood. The "in-sample" pseudo log likelihood is the full log likelihood  $\widehat{\ell}_{in} = \sum_{i=1}^n \widehat{\ell}(\hat{\theta}; \mathfrak{D}_i)$  with the PMLE  $\hat{\theta}$  obtained using the full data, i.e.  $\hat{\theta} = \arg \min \sum_{i=1}^n \widehat{\ell}(\theta; \mathfrak{D}_i)$ . The "out-of-sample" pseudo log likelihood is the cross-validated log likelihood. For  $i = 1, \dots, n$ , let  $\hat{\theta}_{(-i)}$  be the PMLE of  $\theta$  obtained by using the subset of the data which deletes the  $i$ -th observation, i.e.,  $\hat{\theta}_{(-i)} = \arg \min \sum_{j \neq i} \widehat{\ell}(\theta; \mathfrak{D}_j)$ . The "out-of-sample" pseudo log likelihood is given as  $\widehat{\ell}_{out} = \sum_{i=1}^n \widehat{\ell}(\hat{\theta}_{(-i)}; \mathfrak{D}_i)$ . Thus, the PIOS test can be expressed as

$$T_n = \widehat{\ell}_{in} - \widehat{\ell}_{out} \tag{2.4}$$

Intuitively, if the  $i$ -th "in-sample" log-likelihood appears much larger than the "out-of-sample" log-likelihood, then the fitted model will shift to accommodate the  $i$ -th observation, suggesting that the model is in some way inadequate to describe the data. This motivates  $T_n$ , which compares the "in-sample" pseudo log-likelihood and "out-of-sample" pseudo log-likelihood, as a global measure of model adequacy. According to Zhang et al. (2016), under the null hypothesis of correct model specification, the PIOS test statistic  $T_n$  in (2.4) converges in probability to  $p$ , which is the dimension of the parameter vector of  $\theta$ . For example, Archimedean copula family only has one parameter, so  $T_n \xrightarrow{p} 1$  under the null hypothesis. According to Zhang et al. (2016), the PIOS test statistic  $T_n$  is asymptotically equivalent to a so-called information ratio (IR) statistic (Zhou et al., 2012) under the null hypothesis. The IR statistic is defined based on two forms of information matrices, which are negative sensitivity matrix  $S(\theta) \stackrel{\Delta}{=} -\mathbb{E}_0[\ell_{\theta\theta}(\theta; \mathfrak{D}_1)]$  and variability matrix  $V(\theta) \stackrel{\Delta}{=} \mathbb{E}_0[\ell_{\theta}^T(\theta; \mathfrak{D}_1)\ell_{\theta}(\theta; \mathfrak{D}_1)]$ , where  $\ell_{\theta}(\theta; \mathfrak{D}_1) = \frac{\partial}{\partial \theta} l(\theta; \mathfrak{D}_1)$ ,  $\ell_{\theta\theta}(\theta; \mathfrak{D}_1) = \frac{\partial}{\partial \theta \partial \theta^T} l(\theta; \mathfrak{D}_1)$ , and  $\mathbb{E}_0(\cdot)$  is the expectation under the true copula model  $\mathbb{C}_0$ . Under suitable regularity condi-

tions,

$$T_n \xrightarrow{p} \mathbb{E}_0[\ell_{\theta}^T(\theta^*; \mathfrak{D}_1)S^{-1}(\theta)\ell_{\theta^*}(\theta^*; \mathfrak{D}_1)] = \text{tr}\{S^{-1}(\theta^*)V(\theta^*)\}, \text{ as } n \rightarrow \infty,$$

where  $\theta^* \in \Theta$  is the limiting value of  $\widehat{\theta}$ , i.e.,  $\widehat{\theta} \xrightarrow{p} \theta^*$ . Under correct model specification,  $S(\theta^*) = V(\theta^*)$  causes  $\text{tr}\{S^{-1}(\theta^*)V(\theta^*)\} = p$ . The IR statistic is given as  $R_n = \text{tr}\{\widehat{S}^{-1}(\widehat{\theta})\widehat{V}(\widehat{\theta})\}$  where  $\widehat{S}(\theta)$  and  $\widehat{V}(\theta)$  are consistent estimates of the two information matrices  $S(\theta)$  and  $V(\theta)$ , respectively. More proofs can be found in Zhang et al. (2016). In addition, under the null hypothesis,  $R_n$  is asymptotically distributed as a normal random variable. Due to the stochastic equivalence between  $T_n$  and  $R_n$ ,  $T_n$  is also asymptotically normally distributed under the null hypothesis. However, it is difficult to estimate the asymptotic variance analytically. I suggest using bootstrap to approximate the asymptotical variance of the test statistic in finite samples. Specifically, the bootstrap is implemented in the following steps:

Step 1 Sample observations with size  $n$  with replacement

Step 2 Use the re-sampled data to calculate the PIOS test statistic, denoted as  $T_n^{(b)}$

Step 3 Repeat Step 1-2 by  $B$  times

Based on the  $B$  bootstrap counterparts of  $T_n$ , denoted by  $T_n^{\text{B}} = \{T_n^{(b)}, b = 1, \dots, B\}$ , we are able to calculate the standard deviation  $sd\{T_n^{\text{B}}\}$ . The p-value is calculated as  $2 \left[ 1 - \Phi \left( \left| \frac{T_n - 1}{sd\{T_n^{\text{B}}\}} \right| \right) \right]$  where  $\Phi$  is the cdf of a standard normal distribution.

# Chapter 3

## Numerical Illustration

In this chapter, I conduct simulation studies to investigate the finite sample performance of the PIOS test in terms of type I error control and test power. In addition, two real data examples are presented to illustrate the use of PIOS test for data analyses.

### 3.1 Simulation Studies

In this section, I investigate the performance of the PIOS test, such as its empirical type I error and test power, in finite samples. In the simulation studies, I considered four copula models: three Archimedean copulas, which are Frank, Clayton and Gumbel, and the Gaussian copula family. In addition, I investigated the effects of several factors on the performance of the test. These factors include the percentage of censoring, the dependence strength (characterized by Kendall's  $\tau$ ), and the sample size. Specifically, I considered two censoring rates, which are 30% and 60%, three values of Kendall's  $\tau$ , which are 0.2, 0.5, and 0.8, and two sample sizes 500 and 1000. For each dataset, 200 bootstrap samples are used to calculate the empirical  $p$ -value.

#### 3.1.1 Setups

A simulated data set of size  $n$  is generated via the following steps: for each data point  $i = 1, 2, \dots, n$ ,

- Step 1. Simulate a bivariate random vectors, denoted as  $(u_{i1}, u_{i2})$ , from a copula model  $\mathbb{C}_0(\cdot; \cdot; \theta_0)$  with  $\theta_0$  obtained from a given value of Kendall's  $\tau$ .

- 
- Step 2. Generate a bivariate survival times  $(T_{i1}, T_{i2})$  from  $(u_{i1}, u_{i2})$  via  $T_{i1} = \mathbf{F}_1^{-1}(u_{i1})$ ,  $T_{i2} = \mathbf{F}_2^{-1}(u_{i2})$ , where  $\mathbf{F}_j^{-1}(\cdot)$  is the quantile function of a Weibull distribution with rate  $\eta_1 = 2$  and  $\eta_2 = 5$ .
- Step 3. Simulate the censoring time  $C_i$  from an exponential distribution with rate  $d = 0.8$  and  $0.5$ , which results in censoring rate approximately 30% or 60% respectively.
- Step 4. Obtain the observed data  $\mathfrak{D} = \{(X_{i1}, X_{i2}, \delta_{i1}, \delta_{i2}), i = 1, 2, \dots, n\}$ , where  $X_{i1} = \min(T_{i1}, C_i)$ ,  $X_{i2} = \min(T_{i2}, C_i)$ ,  $\delta_{i1} = I(T_{i1} \leq C_i)$ , and  $\delta_{i2} = I(T_{i2} \leq C_i)$ .

Based on the simulated data  $\mathfrak{D} = \{\mathfrak{D}_i, i = 1, \dots, n\}$ , we follow the steps below to implement the PIOS test procedure of testing whether an assumed copula model  $\mathbb{C}(\cdot; \cdot; \theta)$  is correctly specified. The steps are

- Step 1. Obtain the KM estimates of the marginal survival functions for  $T_{i1}$  and  $T_{i2}$  respectively. Obtain the estimates of the density functions by taking the product of the estimated survival functions and estimated hazard functions, where the hazard functions are estimated by the kernel smoothing and interpolation described in Section 2.3. Here, I considered an Epanechnikov kernel function, which is defined as  $K(x) = \frac{3}{4}(1 - x^2)I(|x| \leq 1)$ .
- Step 2. Estimate the PMLE of  $\theta$  by maximizing the pseudo log likelihood function with the estimated marginal survival functions and estimated marginal density functions obtained from the first step. In R, I used `optimize` function for minimizing the negative pseudo log likelihood function by Brent's method. Denote the PMLE by  $\hat{\theta}$ .
- Step 3. Obtain "in-sample" pseudo log-likelihood, which is  $\widehat{\ell}_{in} = \sum_{i=1}^n \widehat{\ell}(\hat{\theta}; \mathfrak{D}_i)$ , where  $\hat{\theta}$  is the PMLE obtained from Step.2.
- Step 4. Obtain "out-of-sample" pseudo log-likelihood, which is  $\widehat{\ell}_{out} = \sum_{i=1}^n \widehat{\ell}(\hat{\theta}_{(-i)}; \mathfrak{D}_i)$ , where  $\hat{\theta}_{(-i)}$  is obtained by maximizing the pseudo log-likelihood with the  $i$ -th data point  $\mathfrak{D}_i$  deleted for  $i = 1, 2, \dots, n$ .
- Step 5. Calculate the PIOS test statistic  $T_n = \widehat{\ell}_{in} - \widehat{\ell}_{out}$ .
- Step 6. Implement the non-parametric bootstrap procedure with  $B = 200$  bootstrap resamples to obtain the empirical  $p$ -value,  $p$ .
- Step 7. Given a significance level  $\alpha$ , if  $p \leq \alpha$ , the null hypothesis is rejected; otherwise, we fail to reject the null hypothesis. Here, I considered the significance level  $\alpha = 0.05$ .

Due to intensive computations required by a large number of simulation scenarios, the results were summarized based on 200 replications for each scenario.

### Type I error control

Table 3.1 and Table 3.2 report the empirical type I errors, which are the empirical proportions of rejecting the null hypothesis among the 200 replications when the null hypothesis is true, at a significance level 5%. The results show that the PIOS test  $T_n$  performed satisfactorily in type I error control for most of the cases. There are some observations I would like to point out. First of all, under each censoring rate and dependence strength, the empirical type I error gets closer to the nominal level as the sample size increases. Especially, when the sample size equals to 1000, the overall performance of the PIOS test is good as expected. Secondly, given the same sample size and censoring rate, type I errors do not demonstrate certain tendencies such as increasing, decreasing or remaining the same value when Kendall's tau increases from 0.2 to 0.8. Thirdly, the empirical type I error decreases as the censoring rate decreases for the sample size 1000 and each value of Kendall's rank correlation. Even though the empirical type I errors are satisfactory in most of the settings, there are some cases with slightly inflated type I errors such under the Clayton model with Kendall's  $\tau = 0.2$ , censor rate = 60% and  $n = 500$ . This problem happens due to an insufficient number of bootstrap resamples.

Table 3.1: Type I error with sample size of 500

True and fitted copula		Gumbel	Frank	Clayton	Gaussian
Censor	Kendall's $\tau$				
0.3	0.2	0.054	0.044	0.065	0.061
	0.5	0.066	0.050	0.058	0.054
	0.8	0.061	0.044	0.060	0.063
0.6	0.2	0.051	0.050	0.067	0.065
	0.5	0.064	0.059	0.064	0.059
	0.8	0.062	0.064	0.058	0.064

### Test Power

Table 3.3 and Table 3.4 report the empirical power of the PIOS test, which is the empirical proportion of correctly rejecting the null hypothesis when the null hypothesis is not true.

Table 3.2: Type I error with sample size of 1000

True and fitted copula		Gumbel	Frank	Clayton	Gaussian
Censor	Kendall's $\tau$				
0.3	0.2	0.052	0.043	0.052	0.056
	0.5	0.041	0.040	0.053	0.042
	0.8	0.046	0.033	0.052	0.044
0.6	0.2	0.038	0.051	0.054	0.058
	0.5	0.038	0.046	0.061	0.050
	0.8	0.044	0.044	0.055	0.056

Here I simulated data from one of the four copulas, and then use the PIOS test for each of the other copulas. Several observations could be drawn from the results in the following

1. In general, the test power increases when the sample size increases, or the dependence strength increases, or the censoring rate decreases.
2. When the Kendall's tau is close to 0, such as  $\tau = 0.2$ , the simulated data are drawn from a copula which is close to the independent copula. Thus, the separation between different copulas become challenging, especially with small sample sizes. It is interesting to note that even in this situation of weak dependence, the PIOS test has demonstrated relatively high power of rejecting the Gaussian copula when the data is simulated from Archimedean copulas with a large sample size.
3. In general, the PIOS test performs well in differentiating between Clayton and Gumbel and between Clayton and Gaussian. It might result from their distinct tail dependence structures where Clayton copulas have lower tail dependence, but Gumbel copulas have upper tail dependence, and Gaussian copulas have no tail dependence, which can be seen in Figure 3.1, which shows the scatter plots of simulated observations from these three copulas with a common Kendall's tau 0.8.
4. There are some cases where the PIOS test performs poorly with low power. For example, when the data are simulated from the Gumbel or Clayton copulas, the power of rejecting the Frank copulas is low, even with large sample size and strong correlation. This also happens when the Gumbel or the Frank copulas are tested but the data simulated from the Gaussian copula. However, if the data are generated from

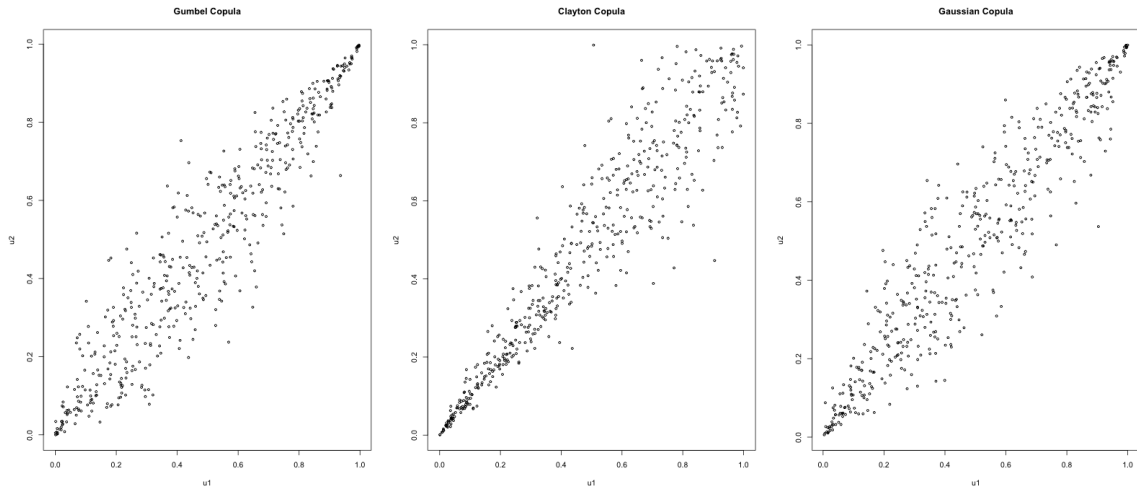


Figure 3.1: Scatter plots of 500 pairs of simulated data from three different copulas with a common Kendall's  $\tau = 0.8$ .

the Gumbel or Frank copulas with moderate correlation, the test will reject the Gaussian copula with high power when the sample size is large and the censoring rate is low. This might be due to the censoring mechanism. The event times are subject to right-censoring, which leads to insufficient information regarding the upper tail dependence from the data. Thus, when the data are generated from the Gumbel copula, it is hard to differentiate the Gumbel copula from the Frank copula which is symmetric among the Archimedean family. Similarly, when the data are simulated from the Gaussian copula, insufficient information on the upper tail dependence will lead to low power of rejecting the Gumbel copula. In contrast, since the Clayton copula has lower tail dependence, it is relatively easier to be differentiated from the other copulas.



Table 3.3: Test power with sample size of 500

True Copula		Gumbel			Frank			Clayton			Gaussian		
Censor	Kendall's $\tau$	Clayton	Frank	Gaussian	Gumbel	Clayton	Gaussian	Gumbel	Frank	Gaussian	Gumbel	Clayton	Frank
0.3	0.2	0.114	0.154	0.403	0.095	0.132	0.200	0.266	0.076	0.240	0.079	0.066	0.056
	0.5	0.882	0.207	0.616	0.150	0.872	0.330	0.854	0.088	0.680	0.174	0.276	0.140
	0.8	0.922	0.297	0.830	0.996	1.000	0.948	0.948	0.372	1.000	0.090	0.858	0.072
0.6	0.2	0.087	0.134	0.412	0.066	0.066	0.182	0.146	0.044	0.216	0.072	0.062	0.053
	0.5	0.750	0.138	0.547	0.120	0.746	0.298	0.684	0.068	0.504	0.146	0.250	0.136
	0.8	0.994	0.114	0.734	0.992	0.992	0.912	1.000	0.250	0.992	0.082	0.784	0.074

Table 3.4: Test power with sample size of 1000

Copula		Gumbel			Frank			Clayton			Gaussian		
Censor	Kendall's $\tau$	Clayton	Frank	Gaussian	Gumbel	Clayton	Gaussian	Gumbel	Frank	Gaussian	Gumbel	Clayton	Frank
0.3	0.2	0.188	0.236	0.754	0.180	0.280	0.584	0.566	0.089	0.640	0.092	0.145	0.062
	0.5	0.992	0.254	0.912	0.420	1.000	0.672	1.000	0.182	0.954	0.230	0.558	0.242
	0.8	1.000	0.430	0.958	1.000	1.000	0.998	1.000	0.845	1.000	0.112	0.990	0.095
0.6	0.2	0.178	0.236	0.772	0.094	0.120	0.516	0.416	0.072	0.576	0.090	0.070	0.058
	0.5	0.974	0.222	0.854	0.402	0.986	0.596	0.964	0.084	0.876	0.218	0.508	0.218
	0.8	1.000	0.271	0.934	1.000	1.000	0.998	1.000	0.628	1.000	0.114	0.978	0.082

---

## 3.2 Application

In this section, I illustrate the PIOS test procedure via two real data examples. The first data example was used in several previous work (Manatunga and Oakes 1999, Wang and Wells 2000, Emura et al. 2010). It is from a diabetic osteopathy study, which examined the effectiveness of laser photocoagulation for delaying the onset of blindness in patients with diabetic retinopathy. According to Manatunga and Oakes (1999), it has been shown that the data set that only includes patients with adult onset diabetes can be fit well using the Clayton model based on the diagnostic plot proposed by Oakes (1989). In addition, Wang (2010) has also shown that the Clayton copula is the best fit to the adult onset data, but different from Manatunga and Oakes (1999), they claimed the Frank copula can also be used to fit the data. To compare our analysis with the previous results, we conducted the PIOS test on the same adult onset data. Approximately 39.8% of them have blindness with treatment, and 78.3% of them experience censoring without the treatment. The empirical estimate of Kendall's rank correlation is 0.376. To conduct the PIOS test, an Epanechnikov kernel function and 2000 bootstrap samples were considered. The corresponding p-value of the PIOS test for the Gumbel, Frank, Clayton, and Gaussian copulas are 0.287, 0.420, 0.483, and 0.119 respectively. At the significance level 0.05, I fail to reject all the four copulas. However, since the p-value for the Clayton copula is the highest among the four, so it seems that the Clayton copula might be the most appropriate model for the adult sub-sample, which agrees with the same conclusion stated in Manatunga and Oakes (1999) and Wang (2010).

The second data example is from the Australian Twin Study (Duffy et al., 1990), in which the ages at appendectomy measured for each twin pairs represents the bivariate event times. A subset of the original data, which contains 748 dizygotic pairs, was studied according to Prentice and Hsu (1997). In this sub-sample, 82 observations were uncensored, 117 were censored for the treated group, 105 were censored for the untreated group, and 444 were censored for both treated and untreated groups. The empirical estimate of Kendall's rank correlation is 0.487. I still considered an Epanechnikov kernel function and 2000 bootstrap samples. The PIOS test was applied on four copula model candidates individually, and the results showed that the Gumbel copula ( $p$ -value = 0.257) and the Frank copula ( $p$ -value = 0.071) are not rejected at the 5% significance level, and the Clayton copula ( $p$ -value = 0.036) and the Gaussian copula ( $p$ -value = 0.000) are rejected at the 5% significance level. I would conclude that the Gumbel copula was the "best" fitting copula to the data, which is consistent with the conclusion drawn by Emura et al. (2010).

# Chapter 4

## Discussion

In this paper, I extended the PIOS test to semi-parametric copula models for right-censored bivariate survival times. The PIOS test statistic is constructed by comparing the "in-sample" pseudo likelihood and "out-of-sample" pseudo likelihood. The essential idea of the PIOS test is to quantify how the assumed copula is sensitive to the change of data. There are several advantages of the PIOS test. First of all, unlike other goodness-of-fit test procedures, our test can be used on all possible copula models that have been considered in the literature. Second, in comparison to the blanket tests considered in Genest et al. (2009) and Scaillet (2007), all of which are rank-based tests, the PIOS test enables us to avoid using any probability integral transformations. Last, the PIOS test is computational straightforward and easily constructed. As shown in simulation studies, it works quite well in type I error control and reach high power in separation of several different copulas. For datasets with small sizes, it still works satisfactorily.

However, there are some issues with the PIOS test. Under the right censoring, some copulas might be difficult to be differentiated from each other. One of the possible reasons might be that the censoring leads to insufficient information on the upper tail dependence. Moreover, the non-parametric bootstrap method considered here might not be able to approximate the finite sample distribution of the test statistic well for some scenarios. Table A.1, A.2, A.3 and A.4 present the average estimated standard errors from bootstrap (boot.se) and the empirical standard errors (em.se). Most of the estimated standard errors are very close to the empirical standard errors. However, when the censoring rate is high, and Kendall's rank correlation is strong, there are some discrepancies. As a consequence, inaccurate estimate of the standard errors might lead to inaccurate power and type I error of the test procedure. Thus, an alternative resampling procedure, such as the parametric bootstrap method, should

---

be considered.

In addition, the assumption of independent censoring is required in the Kaplan-Meier estimates of the survival function. In my future work, I plan to investigate the robustness of the PIOS test against violation of the independent censoring assumption. Moreover, I plan to establish the theoretical proof of the asymptotic properties of the PIOS test in the framework of semi-parametric copula models for right-censored bivariate survival times. In Zhang et al. (2016), the PIOS statistic was shown to be asymptotically equivalent to the IR statistic. Thus, I plan to study the IR statistic for testing the semi-parametric copula models of right-censored event times. Furthermore, to make the computation less time-consuming, I plan to extend the PIOS test by deleting a block of data to construct "out-of-sample" pseudo likelihood. Last but not least, I plan to compare the PIOS test to other existing test procedures to better evaluate its performance.

# Bibliography

- Aalen, O. O. (1978). Nonparametric inference for a family of counting processes. *Annals of Statistics*, 6:701–726.
- Andersen, P. K., Ekstrom, C. T., Klein, J. P., Shu, Y., and Zhang, M.-J. (2005). A class of goodness of fit tests for a copula based on bivariate right-censored data. *Biometrical Journal*, 47:815–824.
- Chen, S. X. and Huang, T.-M. (2007). Nonparametric estimation of copula functions for dependence modelling. *The Canadian Journal of Statistics*, 35:265–282.
- Clayton, D. G. (1978). A model for association in bivariate life tables and its application in epidemiological studies of familial tendency in chronic disease incidence. *Biometrika*, 65:141–151.
- Duffy, D. L., Martin, N. G., and Mathews, J. (1990). Appendectomy in australian twins. *American Journal of Human Genetics*, 47(3):590–592.
- Embrechts, P., McNeil, A., and Straumann, D. (2000). Correlation and dependency in risk management: properties and pitfalls. In *Risk Management*. Cambridge University Press.
- Emura, T., Lin, C.-W., and Wang, W. (2010). A goodness-of-fit test for archimedean copula models in the presence of right censoring. *Computational Statistics & Data Analysis*, 54(12):3033–3043.
- Fermanian, J.-D. (2005). Goodness-of-fit tests for copulas. *Journal of Multivariate Analysis*, 95:119–152.
- Fermanian, J.-D., Radulovic, D., and Wegkamp, M. (2004). Weak convergence of empirical copula processes. *Scandinavian Journal of Statistics*, 10(5):847–860.

- 
- Genest, C., Ghoudi, K., and Rivest, L.-P. (1995). A semiparametric estimation procedure of dependence parameters in multivariate families of distributions. *Biometrika*, 82:543–552.
- Genest, C., Kojadinovic, I., Nešlehová, J., and Jun, Y. (2011). A goodness-of-fit test for bivariate extreme-value copulas. *Bernoulli*, 17:253–275.
- Genest, C., Quessy, J.-F., and Rémillard, B. (2006). Goodness-of-fit procedures for copula models based on the probability integral transformation. *Scandinavian Journal of Statistics*, 33:337–366.
- Genest, C., Rémillard, B., and Beaudoin, D. (2009). Goodness-of-fit tests for copulas: A review and a power study. *Insurance: Mathematics and Economics*, 44:199–213.
- Genest, C. and Segars, J. (2009). Rank-based inference for bivariate extreme-value copulas. *Annals of Statistics*, 37(5B):2990–302.
- Gijbels, I. and Mielniczuk, J. (1990). Estimating the density of a copula function. communications in statistics. *Theory and Methods*, 19:445–464.
- Gribkova, S. and Lopez, O. (2015). Non-parametric copula estimation under bivariate censoring. *Scandinavian Journal of Statistics*, 42:925–946.
- Hougaard, P. (1986). A class of multivariate failure time distributions. *Biometrika*, 73:671–678.
- Huang, W. and Prokhorov, A. (2014). A goodness-of-fit test for copulas. *Econometric Reviews*, 98:533–543.
- Kalbfleisch, J. D. and Prentice, R. L. (2002). *The Statistical Analysis of Failure Time Data. 2nd edition*. John Wiley and Sons.
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, 53(282):457–481.
- Kim, G., Silvapulle, M., and Silvapulle, P. (2007). Comparison of semiparametric and parametric methods for estimating copulas. *Computational Statistics and Data Analysis*, 51(6):2836–2850.
- Li, Y., Prentice, R. L., and Lin, X. (2008). Semiparametric maximum likelihood estimation in normal transformation models for bivariate survival data. *Biometrika*, 95(4):947–960.

- 
- Lin, D. Y. and Ying, Z. (1993). A simple nonparametric estimator of the bivariate survival function under univariate censoring. *Biometrika*, 80:573–581.
- Malevergne, Y. and Sornette, D. (2003). Testing the gaussian copula hypothesis for financial assets dependences. *Quantitative Finance*, 3:231–250.
- Manatunga, A. K. and Oakes, D. (1999). Parametric analysis for matched pair survival data. *Lifetime Data Analysis*, 5:371–387.
- Mesfioui, M., Quessy, J.-F., and Toupin, M.-H. (2009). On a new goodness-of-fit process for families of copulas. *Canadian Journal of Statistics*, 37:80–101.
- Nelson, R. B. (2006). *An Introduction to Copulas*, chapter 4, page 110. Springer.
- Nelson, W. (1972). Theory and applications of hazard plotting for censored failure data. *Technometrics*, 14:945–965.
- Oakes, D. (1989). Bivariate survival models induced by frailties. *Journal of the American Statistical Association*, 84:487–493.
- Omelka, M., Gijbels, I., and Veraverbeke, N. (2009). Improved kernel estimation of copulas: weak convergence and goodness-of-fit testing. *Annals of Statistics*, 37(5B):3023–3058.
- Othus, M. and Li, Y. (2010). A gaussian copula model for multivariate survival data. *Statistics in Biosciences*, 2(2):154–179.
- Prentice, R. L. and Hsu, L. (1997). Regression on hazard ratios and cross ratios in multivariate failure time analysis. *Biometrika*, 35:25–39.
- Prentice, R. L., Moodie, Z. F., and Wu, J. (1997). Nonparametric estimators of the bivariate survival function under simplified censoring conditions. *Biometrika*, 84:863–880.
- Prentice, R. L., Moodie, Z. F., and Wu, J. (2004). Hazard-based nonparametric survivor function estimation. *Journal of the Royal Statistical Society B*, 66:305–319.
- Presnell, B. and Boos, D. (2004). The ios test for model misspecification. *Journal of the American Statistical Association*, 99:216–227.
- Prokhorov, A. and Schmidt, P. (2009). Likelihood-based estimation in a panel setting: robustness, redundancy and validity of copulas. *Journal of Econometrics*, 153:93–104.

- 
- Proschan, F. (1963). Theoretical explanation of observed decreasing failure rate. *Technometrics*, 3(5):375–383.
- Scaillet, O. (2007). Kernel based goodness-of-fit tests for copulas with fixed smoothing parameters. *Journal of Multivariate Analysis*, 98:533–543.
- Shih, J. H. (1998). A goodness-of-fit test for association in a bivariate survival model. *Biometrika*, 85:189–200.
- Shih, J. H. and Louis, T. A. (1995). Inferences on the association parameter in copula models for bivariate survival data. *Biometrics*, pages 1384–1399.
- Sklar, A. (1959). Fonctions de répartition à n dimensions et leurs marges. *Publ. Inst. Stat. Univ. Paris*, 8:229–231.
- Tsai, W. Y., Leurgans, S., and Crowley, J. J. (1997). Nonparametric estimation of a bivariate survival function in presence of censoring. *The Annals of Statistics*, 14:1351–1365.
- Wang, A. (2010). Goodness-of-fit tests for archimedean copula models. *Statistica Sinica*, 20:441–453.
- Wang, J.-L., Muller, H., and Capra, W. (1998). Analysis of oldest-old mortality: Lifetables revisited. *Annals of Statistics*, 26:126–163.
- Wang, W. and Wells, M. T. (2000). Model selection and semiparametric inference for bivariate failure-time data. *Journal of the American Statistical Association*, 95:62–72.
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica*, 50:1–25.
- Zhang, S., Okhrin, O., Zhou, Q., and Song P, P.-K. (2016). Goodness-of-fit test for specification of semiparametric copula dependence models. *Journal of Econometrics*, 193(1):215–233.
- Zhou, Q., Song, P., and Thompson, M. (2012). Information ratio test for model misspecification in quasi-likelihood inference. *Journal of the American Statistical Association*, 107(497):205–213.



# Appendix A

## Empirical standard error comparisons

Table A.1: Standard errors of the Gumbel copula simulation

Sample size of 500

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.085	0.089	0.072	0.071	0.157	0.151	0.126	0.146
	0.5	0.097	0.107	0.094	0.094	0.237	0.217	0.213	0.261
	0.8	0.104	0.114	0.307	0.316	0.405	0.474	0.459	0.542
0.6	0.2	0.093	0.094	0.080	0.077	0.171	0.160	0.139	0.158
	0.5	0.108	0.118	0.107	0.106	0.264	0.258	0.261	0.324
	0.8	0.118	0.128	0.370	0.399	0.456	0.549	0.612	0.717

Sample size of 1000

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.060	0.059	0.050	0.047	0.109	0.105	0.092	0.092
	0.5	0.070	0.068	0.066	0.060	0.174	0.162	0.151	0.163
	0.8	0.073	0.074	0.198	0.183	0.303	0.343	0.288	0.351
0.6	0.2	0.066	0.062	0.055	0.053	0.118	0.117	0.103	0.105
	0.5	0.077	0.076	0.074	0.067	0.195	0.184	0.169	0.184
	0.8	0.083	0.085	0.278	0.272	0.359	0.429	0.363	0.448

Table A.2: Standard errors of the Frank copula simulation

Sample size of 500

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.096	0.095	0.070	0.066	0.154	0.150	0.126	0.142
	0.5	0.125	0.124	0.093	0.095	0.258	0.233	0.196	0.248
	0.8	0.284	0.271	0.231	0.215	0.665	0.680	0.405	0.496
0.6	0.2	0.108	0.103	0.078	0.075	0.167	0.171	0.139	0.154
	0.5	0.137	0.138	0.103	0.109	0.281	0.272	0.211	0.274
	0.8	0.303	0.288	0.261	0.255	0.748	0.817	0.419	0.531

Sample size of 1000

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.067	0.069	0.049	0.047	0.110	0.108	0.094	0.104
	0.5	0.095	0.088	0.065	0.062	0.195	0.176	0.154	0.178
	0.8	0.235	0.221	0.153	0.083	0.525	0.504	0.332	0.385
0.6	0.2	0.073	0.074	0.054	0.052	0.119	0.124	0.105	0.116
	0.5	0.104	0.097	0.073	0.068	0.219	0.199	0.168	0.207
	0.8	0.253	0.239	0.172	0.158	0.603	0.609	0.347	0.406

Table A.3: Standard errors of the Clayton copula simulation

Sample size of 500

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.114	0.094	0.070	0.066	0.122	0.125	0.112	0.120
	0.5	0.137	0.128	0.088	0.083	0.121	0.120	0.185	0.205
	0.8	0.300	0.287	0.193	0.197	0.142	0.175	0.486	0.562
0.6	0.2	0.138	0.104	0.078	0.074	0.132	0.141	0.123	0.136
	0.5	0.145	0.139	0.096	0.094	0.134	0.148	0.191	0.223
	0.8	0.307	0.296	0.184	0.182	0.169	0.200	0.494	0.601

Sample size of 1000

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.072	0.068	0.048	0.048	0.088	0.092	0.084	0.093
	0.5	0.101	0.086	0.062	0.060	0.088	0.098	0.138	0.148
	0.8	0.233	0.203	0.108	0.109	0.096	0.116	0.408	0.436
0.6	0.2	0.079	0.075	0.052	0.051	0.096	0.099	0.091	0.102
	0.5	0.108	0.095	0.068	0.066	0.099	0.112	0.146	0.159
	0.8	0.239	0.211	0.112	0.101	0.117	0.147	0.414	0.452

Table A.4: Standard errors of the Gaussian copula simulation

Sample size of 500

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.084	0.086	0.066	0.063	0.136	0.145	0.104	0.118
	0.5	0.089	0.092	0.077	0.074	0.181	0.199	0.120	0.134
	0.8	0.109	0.128	0.168	0.164	0.251	0.339	0.137	0.159
0.6	0.2	0.095	0.095	0.073	0.071	0.145	0.148	0.114	0.131
	0.5	0.099	0.101	0.085	0.084	0.202	0.232	0.133	0.156
	0.8	0.121	0.140	0.197	0.206	0.290	0.399	0.154	0.190

Sample size of 1000

Censor	Kendall's $\tau$	Gumbel		Frank		Clayton		Gaussian	
		boot.se	em.se	boot.se	em.se	boot.se	em.se	boot.se	em.se
0.3	0.2	0.061	0.061	0.046	0.047	0.097	0.096	0.078	0.081
	0.5	0.064	0.066	0.054	0.055	0.134	0.137	0.086	0.088
	0.8	0.079	0.100	0.085	0.071	0.181	0.219	0.095	0.096
0.6	0.2	0.066	0.067	0.051	0.051	0.106	0.107	0.087	0.093
	0.5	0.071	0.074	0.060	0.060	0.153	0.159	0.097	0.103
	0.8	0.088	0.108	0.117	0.089	0.210	0.268	0.110	0.118