# DYNAMIC COPULA-BASED NONPARAMETRIC ESTIMATION OF RANK-TRACKING PROBABILITIES WITH LONGITUDINAL DATA

Xiaoyu Zhang<sup>1</sup>, Mixia  $Wu^{*1,2}$  and Colin O.  $Wu^3$ 

<sup>1</sup> Beijing University of Technology, <sup>2</sup>Beijing Institute for Scientific and Engineer Computing and <sup>3</sup> National Heart, Lung and Blood Institute

Abstract: The rank-tracking probability (RTP) is a useful statistical index for measuring the "tracking ability" of longitudinal disease risk factors in biomedical studies. A flexible nonparametric method for estimating the RTP is the two-step unstructured kernel smoothing estimator, which can be applied when there are time-invariant and categorical covariates. We propose a dynamic copula-based smoothing method for estimating the RTP, and show that it is both theoretically and practically superior to the unstructured smoothing method. We derive the asymptotic mean squared errors of the copula-based kernel smoothing estimators, and use a simulation study to show that the proposed method has smaller empirical mean squared errors than those of the unstructured smoothing method. We apply the proposed estimation method to a longitudinal epidemiological study and show that it leads to clinically meaningful findings in biomedical applications.

*Key words and phrases:* Dynamic copula model, longitudinal study, rank-tracking probability, risk factor, two-step smoothing, unstructured smoothing.

# 1. Introduction

An important objective of a longitudinal analysis in biomedical studies is to investigate the effect of covariates on the response variables over time. Existing methods focus on regression models based on the conditional mean and correlations. Examples of such methods can be found in Hoover et al. (1998), Rice and Wu (2001), Fan, Huang and Li (2007), and Wu and Tian (2018), among others. These methods, although popular in practice, lack the ability to quantitatively track the persistence of disease risk factors over a time range of interest.

The need for longitudinal methods beyond conditional means and correlations can be demonstrated by the National Heart, Lung, and Blood Institute Growth and Health Study (NGHS), for example, NGHSRG (1992), NHBPEP (2004), and Obarzanek et al. (2010), in which many scientific questions are answered by evaluating the conditional distributions rather than the conditional means or correlations. Designed as a prospective epidemiological study, the

<sup>\*</sup>Corresponding author.

NGHS contains up to 10 annual follow-up observations from 1,213 African American girls and 1,166 Caucasian girls who enrolled in the study at age 9 or 10 years. An important objective of the study is to determine whether a girl's cardiovascular disease risk factor has any tracking ability over an age range.

A class of longitudinal methods in the literature for evaluating tracking is based on modeling serial correlations. However, a serial correlation may be an inadequate measure of tracking ability when the conditional distribution function of the outcome variable is unknown. A practical approach to evaluating tracking ability is to use the conditional distribution functions of the outcome variables, as in Hall, Wolff and Yao (1999) and Hall, Racine and Li (2004). Wu and Tian (2013) studied a class of conditional distributions known as the "rank-tracking probability" (RTP) to quantify the tracking ability of a longitudinal outcome variable, and developed an unstructured kernel smoothing method to estimate the RTPs.

However, the unstructured smoothing estimation proposed by Wu and Tian (2013) is sometimes impractical, because the RTPs involve joint probabilities at two time points that might not be estimated appropriately, owing to a lack of sufficient observations at these time points. This motivates using the copulabased method to estimate the RTPs because, by Sklar's theorem (Sklar (1959)), any multivariate distribution can be expressed by its marginal distributions and a copula function. Based on a given copula model, we can estimate a multivariate distribution by separately estimating the marginal distributions and the copula function. However, if both the marginal distributions and the copula function are estimated using nonparametric estimators, such as the kernel estimators (Scaillet and Fermanian (2002); Chen and Huang (2007)) we may encounter "curse of dimensionality." On the other hand, imposing a parametric structure on both the copula function and the marginal distributions may lead to excessive bias, owing to the potential model misspecification (Härdle et al. (2004)). As a useful compromise, Genest, Ghoudi and Rivest (1995) proposed a flexible semiparametric copula approach in which the copula function is modeled parametrically, but the marginal distributions are estimated nonparametrically by their corresponding empirical distributions. To reduce the potential of a model misspecification, Joe (2014) suggested using a data-driven method that selects the copula function from a set of copula models based on a given model selection criterion.

We develop a dynamic copula-based smoothing method to estimate the RTPs of a longitudinal outcome variable, comprising two estimation steps. First, we estimate the raw joint probabilities at a set of distinct design time points using a copula model selected from a set of candidate copula models by maximizing the likelihood functions. Second, we estimate the dynamic RTPs at any time points by smoothing the raw estimates using a kernel smoothing method (Härdle et al. (2004, Chap. 4)). We compare our copula-based smoothing

method with the unstructured smoothing method of Wu and Tian (2013) using a simulation study. Our simulation results suggest that the proposed method is superior to the unstructured method, because the former has smaller empirical mean squared errors. We apply the proposed estimation method to NGHS blood pressure data, and show that it leads to clinically meaningful results for the tracking patterns of blood pressure levels of adolescent girls.

In Section 2, we introduce the longitudinal data structure and the RTPs. In Section 3, we present the two-step copula-based smoothing estimation procedure, and derive the asymptotic properties of the raw and smoothing estimators. We conduct a simulation study in Section 4, and apply our estimation method to NGHS blood pressure data in Section 5. Section 6 concludes the paper.

#### 2. Data Structure and RTPs

## 2.1. Longitudinal observations at design time points

We consider stochastic processes indexed by the time point  $t \in \mathcal{T}$ , where  $\mathcal{T}$  is a bounded subset of  $[0, \infty)$ . At any given  $t \in \mathcal{T}$ ,  $Y(t) \in R$  is the realvalued outcome variable. For simplicity, our longitudinal sample of  $\{Y(t); t \in \mathcal{T}\}$  is assumed to contain n independent subjects, and each subject is observed at a randomly selected subset of K > 1 distinct "design time points"  $\mathscr{K} =$  $\{t_{(1)}, \ldots, t_{(K)}\}$ , where  $t_{(k)} \in \mathcal{T}$ . For the *i*th subject, for  $1 \leq i \leq n$ , the outcome  $Y_i(t_{ij}) = Y_{ij}$  is collected at time points  $t = t_{ij} \in \mathscr{K} \subset \mathcal{T}$ , for  $j = 1, \ldots, n_i$ , where  $n_i$  is the number of observations for the *i*th subject and  $N = \sum_{i=1}^n n_i$ . At each  $t_{(k)}$ ,  $\mathscr{F}_k$  is the set of subjects with observations when  $\{Y_i(t_{(k)}) : i \in \mathscr{F}_k, k = 1, \ldots, K\}$ is a sample of  $\{Y(t); t \in \mathcal{T}\}$ , where  $Y_i(t_{(k)})$  is the outcome for the *i*th subject. Let  $n_k = \sharp\{i \in \mathscr{F}_k\}$  be the number of subjects in  $\mathscr{F}_k$ , and  $n_{g,h} = \sharp\{i \in \mathscr{F}_g \cap \mathscr{F}_h\}$ be the number of subjects in the intersection of  $\mathscr{F}_g$  and  $\mathscr{F}_h$ .

This formulation of longitudinal samples is common in biomedical studies. In an epidemiological study, a subject's follow-up time is often chosen from a set of "design time points," which may lead to a large K and  $n_{g,h} \ll \min\{n_g, n_h\}$ . When the observed time points are not exactly contained in  $\mathcal{K}$ , it is common to pool adjacent observed time points into  $\mathcal{K}$  using a clinically meaningful criterion.

#### 2.2. The RTPs

Suppose that, for  $t \in \mathcal{T}$ , the health status of a subject at time t is determined by whether  $Y(t) \in A(t)$  for a prespecified subset  $A(\cdot) \subseteq R$ , which may change with t. The tracking ability of  $\{Y(t); t \in \mathcal{K}\}$  at any two time points  $t_1 < t_2$  can be measured using the conditional probability of  $\{Y(t_2) \in A(t_2)\}$  given  $\{Y(t_1) \in A(t_1)\}\}$ , which Wu and Tian (2013) refer to as the RTP based on  $A(\cdot)$  at  $t_1 < t_2$ ,

$$RTP_A(t_1, t_2) = P\{Y(t_2) \in A(t_2) \mid Y(t_1) \in A(t_1)\},$$
(2.1)

where the choice of  $A(\cdot)$  depends on the study questions and scientific objectives. As noted in Wu and Tian (2013), the "rank" in (2.1) does not necessarily refer to a statistical ranking, but is used more generally to characterize the ordinal "health status" of a subject in a biomedical study. A direct extension of (2.1) is to evaluate the probability that the subject's health status develops from status  $A_1(\cdot)$  at time  $t_1$  to status  $A_2(\cdot)$  at time  $t_2$ . The RTP based on  $A_1(\cdot)$  and  $A_2(\cdot)$  at  $t_1 < t_2$  is then

$$RTP_{A_1, A_2}(t_1, t_2) = P\{Y(t_2) \in A_2(t_2) \mid Y(t_1) \in A_1(t_1)\}.$$
(2.2)

In biomedical studies, it is common to define  $A_k(\cdot)$ , for k = 1, 2, using certain threshold values  $y_k(t)$ . A natural choice of  $A_k(t)$  is  $A_k(t) = (y_k(t), \infty)$ , for k = 1, 2, and a threshold-based RTP for (2.2) is

$$RTP_{A_1, A_2}(t_1, t_2) = \frac{P_{A_1, A_2}(t_1, t_2)}{P_{A_1}(t_1)},$$
(2.3)

where  $P_{A_1, A_2}(t_1, t_2) = P\{Y(t_1) > y_1(t_1), Y(t_2) > y_2(t_2)\}$  and  $P_{A_1}(t_1) = P\{Y(t_1) > y_1(t_1)\}$  are the joint and marginal probabilities, respectively, of  $Y(t_2)$  and  $Y(t_1)$ . Furthermore, when  $y_k(t)$  are certain quantile values, Wu and Tian (2013) refer to the corresponding RTP as the "quantile-based RTP." Let  $y_{\alpha_i}(t)$  be the  $(100 \times \alpha_i)$  quantile of Y(t). Then, the quantile-based RTP is

$$RTP_{\alpha_1,\alpha_2}(t_1,t_2) = P\{Y(t_2) > y_{\alpha_2}(t_2) \mid Y(t_1) > y_{\alpha_1}(t_1)\}.$$
(2.4)

Because thresholds are widely used to define the status of a disease, we focus on estimating the threshold-based RTP (2.3) and the quantile-based RTP (2.4).

#### 3. Copula-based Smoothing Estimation

We develop the following two-step estimation procedure for the dynamic RTPs in (2.3) and (2.4) as functions of any two time points  $t_1 < t_2 \in \mathcal{T}$ : (a) obtain the raw estimates for the joint probabilities at the design time points under a copula family; (b) compute the smoothing estimates of the joint probabilities and RTPs at any two time points in  $\mathcal{T}$ .

## 3.1. Dynamic copulas for the joint probabilities

Let  $S_t(y) = P\{Y(t) > y\}$  be the "survival function" of Y(t) at time point t. For any time points  $t_1 < t_2 \in \mathcal{T}$ , the joint "survival function" of  $Y(t_1)$  and  $Y(t_2)$  has the copula expression

$$P(y_1, y_2|t_1, t_2) = P\{Y(t_1) > y_1, Y(t_2) > y_2\} = C_{\theta(t_1, t_2)}(S_{t_1}(y_1), S_{t_2}(y_2)), \quad (3.1)$$

where  $C_{\theta(t_1, t_2)}(s_1, s_2)$  is the "copula function" and  $\theta(t_1, t_2)$  is the unknown timevarying copula parameter at time points  $t_1$  and  $t_2$ . We assume that the copula model of  $C_{\theta(t_1, t_2)}(s_1, s_2)$  is either known or unknown. In the latter case, we assume that it is close to one of the candidate copula models in terms of the Kullback–Leibler distance. This assumption is practical in biomedical studies, because a suitable copula can be elected by evaluating how well the copula model fits the available data. Our objective is to estimate the RTPs as functions of  $(t_1, t_2)$ , defined in (2.3) and (2.4), based on the time-varying copula (3.1) and a smoothing method.

At any given design time points  $t_{(g)} \neq t_{(h)} \in \mathscr{K}$ , by (3.1), we have that  $P_{A_1, A_2}(t_{(g)}, t_{(h)}) = C_{\theta(t_{(g)}, t_{(h)})}(S_{t_{(g)}}(y_1(t_{(g)})), S_{t_{(h)}}(y_2(t_{(h)})))$ , which can be estimated by replacing  $S_{t_{(\cdot)}}(y)$  and  $C_{\theta(t_{(g)}, t_{(h)})}(\cdot, \cdot)$  with their corresponding consistent estimates. For the functional copula parameter  $\theta(t_{(g)}, t_{(h)})$ , we consider two scenarios for its estimation: (a)  $C_{\theta}(s_1, s_2)$  belongs to a known copula model that is specified by the unknown  $\theta = \theta(\cdot)$ ; (b) both the copula model and the functional copula parameter are unknown, but an appropriate model for  $C_{\theta}(s_1, s_2)$ can be selected from a set of candidate copula models.

### 3.2. Estimation with a known copula model

When  $C_{\theta}(s_1, s_2)$  belongs to a known copula model, we first estimate the univariate survival functions  $S_{t_{(\cdot)}}(y)$  using the following empirical survival distribution of  $Y(t_{(j)}) > y$ , for any design time  $t_{(j)}$  and y:

$$\widetilde{S}_{t_{(j)}}(y) = \frac{1}{n_j} \sum_{i \in \mathscr{F}_j} \mathbb{1}_{[Y_i(t_{(j)}) > y]}, \quad j = g, h,$$
(3.2)

where  $1_{[\cdot]}$  is an indicator function. For the copula parameter  $\theta(t_{(g)}, t_{(h)})$ , we assume that  $C_{\theta}(s_1, s_2)$  is differentiable with respect to  $s_1$  and  $s_2$ , and define the derivative  $c_{\theta}(s_1, s_2) = \frac{\partial^2 C_{\theta}(s_1, s_2)}{\partial s_1 \partial s_2}$ , where  $\theta = \theta(t_{(g)}, t_{(h)})$ ,  $s_1 = S_{t_{(g)}}(y(t_{(g)}))$ , and  $s_2 = S_{t_{(h)}}(y(t_{(h)}))$ . From (3.1), we can define the following pseudo log-likelihood function for  $\theta$ :

$$l_{g,h}(\theta \mid C) = \frac{1}{n_{g,h}} \sum_{k \in \left(\mathscr{F}_g \cap \mathscr{F}_h\right)} \log c_\theta \left(\widehat{S}_{kt_{(g)}}, \widehat{S}_{kt_{(h)}}\right), \tag{3.3}$$

where  $\widehat{S}_{kt_{(j)}}$  is the  $n_j/(n_j+1)$  rescaled version of the empirical marginal survival function  $\widetilde{S}_{t_{(j)}}(y)$  at  $y = Y_k(t_{(j)})$  within the set  $\mathscr{F}_j$ , for j = g, h. The rescaling is necessary to avoid the potential unboundedness of  $\log c_{\theta}(s_1, s_2)$  as  $s_1$  or  $s_2$  tend to one. Maximizing the pseudo log-likelihood  $l_{g,h}(\theta \mid C)$  of (3.3) with respect to  $\theta$ , the maximum likelihood estimator of  $\theta(t_{(g)}, t_{(h)})$  is

$$\widehat{\theta}_C(t_{(g)}, t_{(h)}) = \operatorname*{argmax}_{\theta} l_{g,h}(\theta \mid C).$$
(3.4)

A numerical computation of (3.4) can be performed using the procedure described in Genest, Ghoudi and Rivest (1995).

Substituting  $S_{t_{(g)}}(y_1(t_{(g)}))$ ,  $S_{t_{(h)}}(y_2(t_{(h)}))$  and  $\theta(t_{(g)}, t_{(h)})$  in (3.1) with their corresponding estimators (3.2) and (3.4), respectively, we obtain the following raw estimator of  $P_{A_1,A_2}(t_{(g)}, t_{(h)})$  at any design time points  $(t_{(g)}, t_{(h)})$ :

$$\widetilde{P}_{A_1, A_2}\left(t_{(g)}, t_{(h)}\right) = C_{\widehat{\theta}_C(t_{(g)}, t_{(h)})}\left(\widetilde{S}_{t_{(g)}, 1}, \widetilde{S}_{t_{(h)}, 2}\right),$$
(3.5)

where  $\widetilde{S}_{t_{(g)},1} = \widetilde{S}_{t_{(g)}}(y_1(t_{(g)}))$  and  $\widetilde{S}_{t_{(h)},2} = \widetilde{S}_{t_{(h)}}(y_2(t_{(h)}))$ . Specifically, when  $y_1(t_{(g)}) = y_{\alpha_1}(t_{(g)})$  and  $y_2(t_{(h)}) = y_{\alpha_2}(t_{(h)})$ , we have

Specifically, when  $y_1(t_{(g)}) = y_{\alpha_1}(t_{(g)})$  and  $y_2(t_{(h)}) = y_{\alpha_2}(t_{(h)})$ , we have  $S_{t_{(g)}}(y_1(t_{(g)})) = 1 - \alpha_1$  and  $S_{t_{(h)}}(y_2(t_{(h)})) = 1 - \alpha_2$ . Thus, the joint probability  $P_{A_1, A_2}(t_{(g)}, t_{(h)}) = P\{Y(t_{(g)}) > y_{\alpha_1}(t_{(g)}), Y(t_{(h)}) > y_{\alpha_2}(t_{(h)})\}$ , denoted by  $P_{\alpha_1, \alpha_2}(t_{(g)}, t_{(h)})$ , can be estimated by

$$\tilde{P}_{\alpha_1,\,\alpha_2}\left(t_{(g)},t_{(h)}\right) = C_{\hat{\theta}_C(t_{(g)},\,t_{(h)})}(1-\alpha_1,1-\alpha_2),\tag{3.6}$$

where  $\alpha_1$  and  $\alpha_2$  are given in (2.4).

**Remark 1.** Note that  $\widehat{S}_{kt_{(g)}}$ ,  $\widehat{S}_{kt_{(h)}}$ , and  $\widehat{\theta}_C$  are calculated from data sets at different design time points  $\mathscr{F}_g$ ,  $\mathscr{F}_h$ , and  $\mathscr{F}_g \cap \mathscr{F}_h$ , respectively. Thus, the proposed copula estimator of  $P_{A_1, A_2}(t_{(g)}, t_{(h)})$  is less affected by an "unbalanced design" of the data, as in Wu and Tian (2018), and outperforms the unstructured nonparametric estimator

$$\widetilde{P}^{N}_{A_{1},A_{2}}(t_{(g)},t_{(h)}) = \frac{1}{n_{g,h}} \sum_{i \in \left(\mathscr{F}_{g} \cap \mathscr{F}_{h}\right)} \mathbb{1}_{[Y_{i}(t_{(h)}) > y_{1}(t_{(g)}), Y_{i}(t_{(g)}) > y_{2}(t_{(h)})]},$$
(3.7)

which is calculated without using the copula structure (3.1) and relies only on the observations in  $\mathscr{F}_g \cap \mathscr{F}_h$ . The advantage of (3.5) and (3.6) over (3.7) becomes obvious when  $n_{g,h}$  is much smaller than min $\{n_g, n_h\}$ . In addition, for the quantilebased estimation of  $P_{\alpha_1,\alpha_2}(t_{(g)}, t_{(h)})$ , the true percentile curves  $y_{\alpha_1}(t)$  and  $y_{\alpha_2}(t)$ in (3.7) are often unknown and need to be estimated, whereas the copula-based estimator (3.6) does not rely on estimated quantile curves. In the simulation study, we use the true percentile curves for the nonparametric estimators, and we estimate the percentile curves from the corresponding samples in the application to NGHS blood pressure data.

## 3.3. Selection of copula models

In most biomedical studies, the exact copula model  $C_{\theta(t_{(g)}, t_{(h)})}(\cdot, \cdot)$  is unknown, but may be selected from a set of copula models. A reasonable procedure for doing so is to maximize a pseudo-likelihood function among all the candidate copula models; see Joe (2014). If  $\mathcal{M}$  is the set of copula models, then the selected copula model is

$$C^{*}(\cdot) = \operatorname*{argmax}_{C \in \mathcal{M}} l_{g,h}(\widehat{\theta}_{C}|C) = \operatorname*{argmax}_{C \in \mathcal{M}} \left\{ \max_{\theta} l_{g,h}(\theta|C) \right\},$$
(3.8)

where  $l_{g,h}(\theta|C)$  is defined by (3.3). Let  $\widehat{\theta}_{C^*}(t_{(g)}, t_{(h)})$  be the corresponding estimator derived from (3.4) under the selected copula model  $C^*(\cdot)$ . The approximated estimator of  $P_{A_1, A_2}(t_{(g)}, t_{(h)})$  based on  $C^*(\cdot)$  is

$$\widetilde{P}_{A_{1},A_{2}}^{*}\left(t_{(g)},t_{(h)}\right) = C_{\widehat{\theta}_{C^{*}}\left(t_{(g)},t_{(h)}\right)}^{*}\left(\widetilde{S}_{t_{(g)},1},\widetilde{S}_{t_{(h)},2}\right).$$
(3.9)

Furthermore, the corresponding estimator of  $P_{\alpha_1,\alpha_2}(t_{(g)},t_{(h)})$  is

$$\widetilde{P}^*_{\alpha_1,\,\alpha_2}(t_{(g)},t_{(h)}) = C^*_{\widehat{\theta}_{C^*}(t_{(g)},\,t_{(h)})}(1-\alpha_1,1-\alpha_2)\,.$$
(3.10)

Because the true joint probabilities may not necessarily belong to the selected copula model, we refer to (3.10) as an "approximated estimator", because the selected copula model may only be a reasonable approximation of the true copula model. The following remarks clarify some implications of relying on an approximated copula model in (3.10).

**Remark 2.** It follows from Joe (2014) that if the Kullback–Leibler distance between the true and selected copula model is small, then so is the difference between the joint probabilities derived from the true copula and the selected copula. Thus, a good choice of candidate copula models  $\mathcal{M}$  can lead to an appropriate estimator  $\tilde{P}^*_{A_1,A_2}(t_{(g)},t_{(h)})$  or  $\tilde{P}^*_{\alpha_1,\alpha_2}(t_{(g)},t_{(h)})$  that is close to  $\tilde{P}_{A_1,A_2}(t_{(g)},t_{(h)})$  or  $\tilde{P}_{\alpha_1,\alpha_2}(t_{(g)},t_{(h)})$ , respectively, under the true copula.

**Remark 3.** As discussed in Joe (2014), the similarity of copulas depends on the closeness of the dependence in the tails. This suggests that the tail properties are useful for identifying distribution functions in copula model selection. In particular, the Frank copula is symmetric and has no tail dependence, whereas the Clayton and Gumbel copulas have strong lower and upper tail dependence, respectively. Because these three copula models capture most of the dependence structures seen in real applications, they are widely used as candidate copula models in the literature. We demonstrate in our simulation study, discussed in Section 4, that the estimators based on the selected copula from these three candidate copula models give satisfactory performance in practice.

## 3.4. Smoothing estimators

For the smoothing step of the procedure, we use the raw estimates  $\tilde{P}_{A_1, A_2}(t_{(g)}, t_{(h)})$  at  $t_{(g)} \neq t_{(h)} \in \mathcal{K}$  in (3.9) to estimate  $P_{A_1, A_2}(t_1, t_2)$  at any time points  $t_1 < t_2 \in \mathcal{T}$  using a kernel smoothing method. The smoothing estimator of  $P_{A_1, A_2}(t_1, t_2)$  is then given by

$$\widehat{P}_{A_1,A_2}(t_1,t_2) = \sum_{g \neq h}^{J} W_{g,h}(t_1,t_2) \widetilde{P}_{A_1,A_2}(t_{(g)},t_{(h)}), \qquad (3.11)$$

where  $W_{g,h}(t_1, t_2)$  is kernel-based weight function,

$$W_{g,h}(t_1, t_2) = \frac{n_{g,h} K\big((t_1 - t_{(g)})/h_1, (t_2 - t_{(h)})/h_2\big)}{\sum_{g \neq h} n_{g,h} K\big((t_1 - t_{(g)})/h_1, (t_2 - t_{(h)})/h_2)\big)},$$
(3.12)

 $K(\cdot, \cdot)$  is a bivariate nonnegative kernel function, and  $h_1$  and  $h_2$  are the corresponding bandwidths.

Similarly, the kernel smoothing estimator of  $P_{A_1,A_2}(t_1,t_2)$  based on the selected copula  $C^*(\cdot)$  is given by

$$\widehat{P}^*_{A_1, A_2}(t_1, t_2) = \sum_{g \neq h}^J W_{g, h}(t_1, t_2) \widetilde{P}^*_{A_1, A_2}(t_{(g)}, t_{(h)}).$$
(3.13)

The smoothing estimator of the marginal probability  $P_{A_1}(t_1) = S_{t_1}(y_1(t_1))$ , based on  $\widetilde{S}_{t_{(k)}}(y_1(t_{(k)}))$ , for  $k = 1, \ldots, K$ , is given by

$$\widehat{S}_{t_1}(y_1(t_1)) = \sum_{k=1}^{K} W_k(t_1) \widetilde{S}_{t_{(k)}}(y_1(t_{(k)})), \qquad (3.14)$$

where  $W_k(t_1) = n_k K((t_1 - t_{(k)})/h_0) / \sum_{j=1}^K n_j K((t_1 - t_{(j)})/h_0), K(\cdot)$  is a nonnegative kernel function, and  $h_0$  is the corresponding bandwidth.

It follows from (3.11) and (3.14) that the kernel smoothing estimator of  $RTP_{A_1,A_2}(t_1,t_2)$  in (2.3) is

$$\widehat{RTP}_{A_1, A_2}(t_1, t_2) = \frac{\widehat{P}_{A_1, A_2}(t_1, t_2)}{\widehat{S}_{t_1}(y_1(t_1))}$$
(3.15)

when the true copula model  $C(\cdot)$  is known, and

$$\widehat{RTP}^*_{A_1, A_2}(t_1, t_2) = \frac{\widehat{P}^*_{A_1, A_2}(t_1, t_2)}{\widehat{S}_{t_1}(y_1(t_1))}$$
(3.16)

when the copula model  $C^*(\cdot)$  is selected from the set of copula models  $\mathcal{M}$ . For the quantile-based RTP (2.4), because  $P_{\alpha_i}(t_i) = P(Y(t_i) > y_{\alpha_i}(t_i)) = 1 - \alpha_i$  for each  $t_i \in \mathcal{T}$ , for i = 1, 2, the kernel smoothing estimators (3.15) and (3.16) can be simplified as

$$\widehat{RTP}_{\alpha_1, \alpha_2}(t_1, t_2) = \frac{\widehat{P}_{\alpha_1, \alpha_2}(t_1, t_2)}{1 - \alpha_i}$$
(3.17)

and

$$\widehat{RTP}^{*}_{\alpha_{1},\alpha_{2}}(t_{1},t_{2}) = \frac{\widehat{P}^{*}_{\alpha_{1},\alpha_{2}}(t_{1},t_{2})}{1-\alpha_{i}},$$
(3.18)

896

respectively, where

$$\widehat{P}_{\alpha_1,\,\alpha_2}(t_1,t_2) = \sum_{g \neq h}^{J} W_{g,\,h}(t_1,t_2) \widetilde{P}_{\alpha_1,\,\alpha_2}(t_{(g)},t_{(h)}),\tag{3.19}$$

and

$$\widehat{P}^*_{\alpha_1,\,\alpha_2}(t_1,t_2) = \sum_{g \neq h}^J W_{g,\,h}(t_1,t_2) \widetilde{P}^*_{\alpha_1,\,\alpha_2}(t_{(g)},t_{(h)}).$$
(3.20)

**Remark 4.** It is well known in the literature that, in practice, the appropriateness of kernel smoothing estimators is affected most by the bandwidth choices, whereas the kernel functions are relatively less important. Thus, several kernel functions may be used to compute (3.17) and (3.18). For simplicity, we use the product kernels  $K(u_1, u_2) = K_1(u_1)K_2(u_2)$  with the same bandwidth  $h_1 = h_2 = h$  in all numerical computations in this paper. Specifically, we use the Gaussian kernel  $K_1(u_1) = (2\pi)^{-1/2} \exp\{-(u_1^2/2)\}$  for  $\hat{P}_{\alpha_1,\alpha_2}(t_1,t_2)$ . For data-driven bandwidth choices of h, we use the "leave-one-subject-out cross-validation" (LSCV) approach described in Wu and Tian (2018, Sec. 12.3.5).

#### 3.5. Asymptotic Properties

In this section, we present the consistency and asymptotic normality of the raw estimators, the consistency of the dynamic copula function estimators, and the asymptotic mean squared risks of the kernel smoothing RTP estimators.

## 3.5.1. Asymptotic properties of raw probability estimators

We first consider the asymptotic properties of the raw estimators (3.5) and (3.6) when the copula model is known. Assume that the joint survival distribution function of  $Y(t_1)$  and  $Y(t_2)$ , for any  $t_1 < t_2 \in \mathcal{T}$ , belongs to a known copula model that satisfies the following assumptions:

- C1. For l = 1, 2, the partial derivatives  $\partial \log c_{\theta}(s_1, s_2)/\partial \theta$ ,  $\partial \log c_{\theta}(s_1, s_2)/\partial s_l$ ,  $\partial^2 \log c_{\theta}(s_1, s_2)/\partial \theta^2$ , and  $\partial^2 \log c_{\theta}(s_1, s_2)/\partial \theta \partial s_l$  are all continuous and can be bounded by  $dC_{\theta}$  integrable functions for any  $(s_1, s_2) \in (0, 1)^2$ .
- C2. For any fixed  $t_1 < t_2 \in \mathcal{T}$ , the functional Fisher information  $I(\theta|t_1, t_2) = -E_{\theta} \{\partial \log c_{\theta} [S_{t_1}(Y(t_1)), S_{t_2}(Y(t_2))]/\partial \theta\}^2$  is finite and bounded away from zero, that is,  $0 < I(\theta|t_1, t_2) < \infty$ .

The following theorem shows that the raw estimator  $\widetilde{P}_{A_1,A_2}(t_{(g)},t_{(h)})$  of (3.5) is consistent and that  $\widetilde{P}_{\alpha_1,\alpha_2}(t_{(g)},t_{(h)})$  of (3.6) is asymptotically normal when there are many subjects with observations at the design time points  $(t_{(g)},t_{(h)})$ .

**Theorem 1.** If C1–C2 hold, then, for any two design time points  $t_{(g)} < t_{(h)} \in \mathcal{K}$ , we have that, as  $n_{g,h} \to \infty$ ,  $\tilde{P}_{A_1,A_2}(t_{(g)},t_{(h)}) \xrightarrow{P} P_{A_1,A_2}(t_{(g)},t_{(h)})$  and  $\sqrt{n_{g,h}}$ 

 $\{\widetilde{P}_{\alpha_1,\,\alpha_2}(t_{(g)},t_{(h)}) - P_{\alpha_1,\,\alpha_2}(t_{(g)},t_{(h)})\} \xrightarrow{L} N(0,\sigma^2(t_{(g)},t_{(h)})), \text{ where } \sigma^2(t_{(g)},t_{(h)}) \text{ is the asymptotical variance.}$ 

A proof for this theorem and an explicit expression of  $\sigma_{g,h}$  are given in the online Supplemental Material.

## 3.5.2. Consistency of dynamic copula estimators

When the copula model  $C(\cdot)$  is unknown, we show that, under the best copula model  $C^*(\cdot)$  from the collection of copula models  $\mathcal{M}$ , the approximated raw estimator (3.10) is consistent when there are many subjects with observations at the design time points  $(t_{(g)}, t_{(h)})$ .

First, by the derivations in Joe (2014), we have the following lemma, which shows that, when the selected copula  $C^*(\cdot)$  is "close" to the true copula  $C(\cdot)$ , the estimated functional copula parameter converges to the functional copula parameter under  $C^*(\cdot)$ .

**Lemma 1.** If  $c_{\theta}(s_1, s_2)$  and  $c_{\theta^*}^*(s_1, s_2)$  are the densities of the true copula  $C(\cdot)$ and the selected copula  $C^*(\cdot)$ , respectively, at time points  $(t_{(g)}, t_{(h)})$ , then  $\hat{\theta}_{C^*}(t_{(g)}, t_{(h)}) \xrightarrow{P} \theta_{C^*}(t_{(g)}, t_{(h)})$  in probability as  $n \to \infty$ , where  $\theta_{C^*}(t_{(g)}, t_{(h)}) = \operatorname{argmax}_{\theta^*} \iint c_{\theta}(s_1, s_2) \log c_{\theta^*}^*(s_1, s_2) \operatorname{ds}_1 \operatorname{ds}_2$ ; that is,  $\theta_{C^*}(t_{(g)}, t_{(h)})$  is the copula parameter such that  $c_{\theta_{C^*}(t_{(g)}, t_{(h)})}^*(s_1, s_2)$  is the closest copula density to  $c_{\theta}(s_1, s_2)$  among all copula densities of the form  $c_{\theta^*}^*(s_1, s_2)$  in the Kullback–Leibler divergence.

Let  $P_{A_1,A_2}^*(t_{(g)},t_{(h)}) = C_{\theta_{C^*}}^*(S_{t_{(g)}}(y_1(t_{(g)})), S_{t_{(h)}}(y_2(t_{(h)})))$  and  $P_{\alpha_1,\alpha_2}^*(t_{(g)},t_{(h)})$ =  $C_{\theta_{C^*}(t_{(g)},t_{(h)})}^*(1-\alpha_1,1-\alpha_2)$ . Using the continuous mapping theorem, the following theorem shows that, under the selected copula model, the raw approximated estimator (3.10) is a consistent estimator of  $P_{\alpha_1,\alpha_2}^*(t_{(g)},t_{(h)})$  at the design time points  $(t_{(g)},t_{(h)})$ .

**Theorem 2.** If the selected copula model  $C^*(\cdot)$  satisfies the assumptions C1–C2, then, for any two design time points  $t_{(g)} < t_{(h)} \in \mathscr{K}$ , as  $n_{g,h} \to \infty$ ,  $\widetilde{P}^*_{A_1,A_2}(t_{(g)}, t_{(h)}) \xrightarrow{P} P^*_{A_1,A_2}(t_{(g)}, t_{(h)})$  and

$$\sqrt{n_{g,h}} \left\{ \widetilde{P}^*_{\alpha_1,\alpha_2}(t_{(g)}, t_{(h)}) - P^*_{\alpha_1,\alpha_2}(t_{(g)}, t_{(h)}) \right\} \xrightarrow{L} N(0, \sigma^2(t_{(g)}, t_{(h)})).$$

This asymptotic result suggests that even if the true copula model is unknown, the raw approximated estimator (3.9) ( or (3.10)) could still be close to the target probability  $P_{A_1,A_2}^*(t_{(g)},t_{(h)})$  (or  $P_{\alpha_1,\alpha_2}^*(t_{(g)},t_{(h)})$ ) when the number of observations in  $\mathscr{K}$  is large.

## 3.5.3. Asymptotic risk of the smoothing RTP estimators

We now derive the asymptotic mean squared errors of the kernel smoothing RTP estimators (3.17) and (3.18), under the following assumptions:

C3. For all  $t \in \mathcal{T}$ ,  $S_t(y)$  is twice continuously differentiable with respect to t.

- C4. For all  $t_1, t_2 \in \mathcal{T}$ , the positive density function  $f(t_1, t_2)$  is continuously differentiable with respect to  $t_1$  and  $t_2$ .
- C5. Let  $\boldsymbol{u} = (u_1, u_2)$ . The bivariate kernel  $K(u_1, u_2) = K_1(u_1)K_2(u_2)$  is a symmetric probability density on a bounded set, assumed to be  $[-1, 1]^2$ , and satisfies  $\int \boldsymbol{u}K(\boldsymbol{u})d\boldsymbol{u} = \boldsymbol{0}$ ,  $R(K) = \int K^2(\boldsymbol{u})d\boldsymbol{u}$ ,  $\int u_1^2 K_1(u_1)du_1 = \mu_{(21)}$  and  $\int u_2^2 K_2(u_2)du_2 = \mu_{(22)}$ , for some positive constants  $\mu_{(21)}$  and  $\mu_{(22)}$ .
- C6. h satisfies that  $h \to 0$  and  $nh^2 \to \infty$  as  $n \to \infty$ .
- C7. For all i = 1, ..., n and  $j_1 \neq j_2, |t_{ij_1} t_{ij_2}| > h$ .

The following theorem shows the consistency of the proposed estimator (3.15), and provides asymptotic expressions for the bias and variance of the estimator (3.17).

**Theorem 3.** If the number of subjects n is large,  $t_1 < t_2$  are interior points within the support of  $f(\cdot, \cdot)$ , and assumptions C1–C7 are satisfied, then

$$\widehat{RTP}_{A_1,A_2}(t_1,t_2) \xrightarrow{P} RTP_{A_1,A_2}(t_1,t_2) \quad as \ n \to \infty.$$
(3.21)

Specifically,  $\widehat{RTP}_{\alpha_1,\alpha_2}(t_1,t_2) \xrightarrow{P} RTP_{\alpha_1,\alpha_2}(t_1,t_2)$  as  $n \to \infty$ , and the asymptotic bias and variance of  $\widehat{RTP}_{\alpha_1,\alpha_2}(t_1,t_2)$  in (3.17) are, respectively,

$$\begin{split} Bias \Big[ \widehat{RTP}_{\alpha_1, \, \alpha_2}(t_1, t_2) \Big] \\ &= \Big[ \mu_{(21)} \left( \frac{\partial P_{\alpha_1, \, \alpha_2}(t_1, t_2)}{\partial t_1} \frac{\partial f(t_1, t_2)}{\partial t_1} + \frac{1}{2} \frac{\partial^2 P_{\alpha_1, \, \alpha_2}(t_1, t_2)}{\partial t_1^2} f(t_1, t_2) \right) \\ &+ \mu_{(22)} \left( \frac{\partial P_{\alpha_1, \, \alpha_2}(t_1, t_2)}{\partial t_2} \frac{\partial f(t_1, t_2)}{\partial t_2} + \frac{1}{2} \frac{\partial^2 P_{\alpha_1, \, \alpha_2}(t_1, t_2)}{\partial t_2^2} f(t_1, t_2) \right) \Big] \\ &\times (1 - \alpha_1)^{-1} f(t_1, t_2)^{-1} h^2 + o(h^2) \end{split}$$

and

$$\begin{split} Var \Big[ \widehat{RTP}_{\alpha_1,\alpha_2}(t_1,t_2) \Big] &= (1-\alpha_1)^{-2} \left( N \, h^2 f(t_1,t_2) \right)^{-1} \sigma^2(t_1,t_2) R(K) \\ &+ o \bigg( \frac{1}{N \, h^2} \bigg), \end{split}$$

where  $R(K) = \int K^2(\boldsymbol{u}) d\boldsymbol{u}$ .

The next theorem shows the convergence of the proposed estimator (3.16), and provides asymptotic expressions for the bias and variance of the proposed estimator (3.18) under the selected copula model.

**Theorem 4.** If the selected copula model  $C^*(\cdot)$  satisfies assumptions C1–C2 and assumptions C3–C7 hold, then

$$\widehat{RTP}^*_{A_1,A_2}(t_1,t_2) \xrightarrow{P} RTP^*_{A_1,A_2}(t_1,t_2) \text{ as } n \to \infty.$$

Specifically,  $\widehat{RTP}^*_{\alpha_1,\alpha_2}(t_1,t_2) \xrightarrow{P} RTP^*_{\alpha_1,\alpha_2}(t_1,t_2)$  as  $n \to \infty$ , and the asymptotic bias and variance of  $\widehat{RTP}^*_{\alpha_1,\alpha_2}(t_1,t_2)$  in (3.18) are, respectively,

$$\begin{split} Bias \left[ \widehat{RTP}^{*}_{\alpha_{1},\,\alpha_{2}}(t_{1},t_{2}) \right] \\ &= \left[ \mu_{(21)} \left( \frac{\partial P^{*}_{\alpha_{1},\,\alpha_{2}}(t_{1},t_{2})}{\partial t_{1}} \frac{\partial f(t_{1},t_{2})}{\partial t_{1}} + \frac{1}{2} \frac{\partial^{2} P^{*}_{\alpha_{1},\,\alpha_{2}}(t_{1},t_{2})}{\partial t_{1}^{2}} f(t_{1},t_{2}) \right) \\ &+ \mu_{(22)} \left( \frac{\partial P^{*}_{\alpha_{1},\,\alpha_{2}}(t_{1},t_{2})}{\partial t_{2}} \frac{\partial f(t_{1},t_{2})}{\partial t_{2}} + \frac{1}{2} \frac{\partial^{2} P^{*}_{\alpha_{1},\,\alpha_{2}}(t_{1},t_{2})}{\partial t_{2}^{2}} f(t_{1},t_{2}) \right) \right] \\ &\times (1-\alpha_{1})^{-1} f(t_{1},t_{2})^{-1} h^{2} + o(h^{2}) + (1-\alpha_{1})^{-1} b(t_{1},t_{2}) \end{split}$$

and  $Var[\widehat{RTP}^*_{\alpha_1,\alpha_2}(t_1,t_2)] = (1 - \alpha_1)^{-2} \{Nh^2 f(t_1,t_2)\}^{-1} \sigma^2(t_1,t_2) R(K) + o(1/(Nh^2)), \text{ where } b(t_1,t_2) = P^*_{\alpha_1,\alpha_2}(t_1,t_2) - P_{\alpha_1,\alpha_2}(t_1,t_2).$ 

**Remark 5.** It follows from Theorems 3 and 4 that if the true copula model is known or  $P_{A_1,A_2}^*(t_1,t_2) = P_{A_1,A_2}(t_1,t_2)$ , then the copula-based kernel smoothing estimator (3.15) or (3.16) of the RTP is asymptotically unbiased. In fact, the proposed estimator (3.16) can be close to  $RTP_{A_1,A_2}(t_1,t_2)$  only if the Kullback-Leibler distance between the true copula model and one of the candidate copula models is small. In addition,  $RTP_{\alpha_1,\alpha_2}(t_1,t_2)$  and  $RTP_{\alpha_1,\alpha_2}(t_1,t_2)$  have the same asymptotic variance, whereas the asymptotic biases of (3.17) and (3.18) differ by a fraction of  $b(t_1,t_2)$  that depends on the difference between  $P_{\alpha_1,\alpha_2}^*(t_1,t_2)$ and  $P_{\alpha_1,\alpha_2}(t_1,t_2)$ .

## 4. Simulation Study

We conduct a simulation study to investigate the finite-sample properties of the proposed copula-based smoothing estimators, and then compare them with those of the unstructured nonparametric smoothing estimators proposed by Wu and Tian (2013). Let  $\{t_{(1)}, \ldots, t_{(40)}\} = \{0.25, 0.5, \ldots, 10\}$  be the design time points. We generate 1,000 subjects with 10 visits per subject in each sample. The *j*th visit time of the *i*th subject  $t_{ij}$  corresponds to  $t_{(k_{ij})}$  in the set of design time points. For each subject  $i = 1, \ldots, 1000$  at the *j*th time point  $t_{ij}$ , we generate the observation  $Y_{ij} = Y_i(t_{ij})$  from

$$Y_{ij} = 21.5 + 0.7(t_{ij} - 5) - 0.05(t_{ij} - 5)^2 + \epsilon_{ij}, \quad j = 1, \dots, 40,$$
(4.1)

where  $\epsilon_{ij} = z_{ik_{ij}}$  is the  $k_{ij}$ th element of  $(z_{i1}, \ldots, z_{i40})$ , which are generated from the 40-dimensional *t*-copula  $C^t_{\mathbf{R}, v}(u_1, \ldots, u_{40})$  with the dispersion structure

Table 1. Empirical biases and RMSEs of the copula-based smoothing estimates and the unstructured nonparametric smoothing estimates obtained from 1,000 simulation replications based on  $A_1(t_1) = (y_{\alpha_1}(t_1), \infty)$  and  $A_2(t_2) = (y_{\alpha_2}(t_2), \infty)$ , with  $\alpha_1 = \alpha_2 = 0.8$ .

RTP	$t_2$	True	Copula Smoothing		Unstructured Smoothing	
			Bias	RMSE	Bias	RMSE
$RTP_{0.8, 0.8}(t_2 - 3, t_2)$	5	0.352	0.003	0.014	0.009	0.038
	6	0.352	0.003	0.014	0.008	0.038
	7	0.352	0.003	0.015	0.009	0.039
	8	0.352	0.003	0.015	0.007	0.040
	9	0.352	0.010	0.018	0.007	0.040
$RTP_{0.8, 0.8}(3, t_2)$	5	0.424	-0.005	0.015	0.011	0.043
	6	0.352	0.003	0.014	0.008	0.038
	7	0.309	-0.001	0.014	0.004	0.035
	8	0.283	-0.005	0.015	0.003	0.036
	9	0.266	-0.006	0.015	0.003	0.034

$$\mathbf{R} = \begin{pmatrix} 1 & \rho & \cdots & \rho^{39} \\ \rho & 1 & \cdots & \rho^{38} \\ \vdots & \vdots & & \vdots \\ \rho^{39} & \rho^{38} & \cdots & 1 \end{pmatrix},$$

where  $\rho = 0.9$ , and v = 4 degrees of freedom, and  $k_{ij}$  is generated by the ceiling of a random number from the uniform distribution U[4(j-1), 4j], for  $j = 1, 2, \ldots, 10$ .

We consider the three most commonly used Archimedean copulas, namely, the Frank, Clayton, and Gumbel copulas, as our candidate copula models. In our simulation, the copula model used in the smoothing estimators is not necessarily from the true t-copula model, but is an Archimedean copula that is "closest" to the true t-copula from which the data are generated. The simulation has 1,000 replications.

We first consider the quantile-based RTPs. Let  $A_1(t_1) = (y_{\alpha_1}(t_1), \infty)$  and  $A_2(t_2) = (y_{\alpha_2}(t_2), \infty)$ , with  $\alpha_1 = \alpha_2 = 0.8$ . We calculate the empirical bias and root mean squared error (RMSE) for each estimator (3.18) of  $RTP_{0.8,0.8}(t_2 - 3, t_2)$  and  $RTP_{0.8,0.8}(3, t_2)$  at a sequence of  $t_2$  values,  $t_1 = t_2 - 3$  and  $t_1 = 3$ . Here,  $RTP_{0.8,0.8}(t_2 - 3, t_2)$  and  $RTP_{0.8,0.8}(3, t_2)$  represent the "three-year tracking ability" and the "first  $(t_2 - 3)$ -year tracking ability," respectively, for the simulated samples. For comparison, we also present the true RTP values and compute the empirical biases and RMSEs from the unstructured smoothing estimates.

Table 1 shows the true RTP values and the empirical biases and RMSEs obtained from the copula-based and unstructured smoothing estimates at several selected time points. Table 1 shows that both smoothing estimates have small



Figure 1. The true RTP curves, the averages of the estimated  $RTP_{\alpha_1, \alpha_2}(t_1, t_2)$  with  $\alpha_1 = \alpha_2 = 0.8$  using the copula-based smoothing estimator and the unstructured smoothing estimator, and the lower and upper 2.5th percentiles computed from 1,000 simulated samples generated from (4.1).

biases. However, the copula-based smoothing estimates are superior to the unstructured smoothing estimates in the sense that they have smaller RMSEs for all the time points shown in Table 1. This superiority holds even when the candidate Archimedean copula models do not include the true t-copula model, and when the unstructured smoothing estimates are calculated using the true percentile curve  $y_{0.8}(t)$ .

Figure 1 shows the averages of the estimated  $RTP_{0.8,0.8}(t_2 - 3, t_2)$  in Figure 1(a) and  $RTP_{0.8,0.8}(3, t_2)$  in Figure 1(b), as well as their corresponding lower and upper 2.5th percentiles computed from 1,000 simulated samples using the copula-based smoothing estimator and the unstructured smoothing estimator. The figure shows that the averages of the estimated curves from both estimators are close to the true RTP curves. However, the gaps between the lower and upper 2.5th percentile curves of the copula-based smoothing estimates are much smaller than those of the unstructured smoothing estimates, suggesting that the copula-based smoothing estimator exhibits less variability than does the unstructured smoothing estimator. This result is consistent with the empirical RMSE results in Table 1.

We next consider the threshold-based RTPs. Let  $A_1(t_1) = (28, \infty)$  and  $A_2(t_2) = (28, \infty)$ . The corresponding empirical biases, RMSEs, averages, and percentiles of the estimator (3.16) and the unstructured smoothing estimator are summarized in Table 2 and Figure 2. These results again suggest that the copula-based estimator (3.16) leads to similar empirical biases, but smaller RMSEs and less variability compared with those of the corresponding unstructured smoothing estimator.

Table 2. Empirical biases and RMSEs of the copula-based smoothing estimates and the unstructured nonparametric smoothing estimates obtained from 1,000 simulation replications based on  $A_1(t_1) = (28, \infty)$  and  $A_2(t_2) = (28, \infty)$ .

RTP	$t_2$	True	Copula Smoothing		Unstructured Smoothing	
			Bias	RMSE	Bias	RMSE
$RTP_{A_1, A_2}(t_2 - 3, t_2)$	5	0.430	-0.005	0.032	0.005	0.043
	6	0.442	-0.001	0.031	0.005	0.039
	7	0.452	0.001	0.029	0.002	0.038
	8	0.460	0.003	0.028	0.003	0.035
	9	0.466	0.004	0.027	0.005	0.033
$RTP_{A_1, A_2}(3, t_2)$	5	0.495	-0.001	0.032	0.010	0.041
	6	0.442	-0.001	0.031	0.005	0.039
	7	0.412	-0.001	0.031	0.003	0.038
	8	0.396	-0.004	0.032	0.003	0.039
	9	0.388	-0.006	0.032	0.001	0.037



Figure 2. The true RTP curves, the averages of the estimated  $RTP_{A_1, A_2}(t_1, t_2)$  with  $A_1(t) = (28, \infty), A_2(t) = (28, \infty)$  using the copula-based smoothing estimator and the unstructured smoothing estimator, and the lower and upper 2.5th percentiles computed from 1,000 simulated samples generated from (4.1).

Tables 1 and 2 and Figures 1 and 2 demonstrate that the copula-based smoothing estimator is superior to the unstructured smoothing estimator, both when estimating quantile-based RTPs and when estimating threshold-based RTPs. The only difference between these two scenarios is that the marginal probabilities  $S(t_1)$  and  $S(t_2)$  are equal to  $1 - \alpha_1$  and  $1 - \alpha_2$ , respectively, in (3.18), whereas they need to be estimated in (3.16).

## 5. Application to NGHS Blood Pressure Data

We apply our estimation method to NGHS blood pressure (BP) data. This data set is described and analyzed by Wu and Tian (2013) and Wu and Tian (2018, Chap. 12) using the unstructured nonparametric smoothing estimator. Given that an important objective of the NGHS is to evaluate the racial differences of BP distributions, we estimate the RTPs for Caucasian girls and African American girls separately.

Because age (in years) in the data set is rounded up to six decimal places, there are many distinct time points, but few subjects observed at each distinct time point. A practical approach that is clinically meaningful and has been used in the literature is "data binning," which pools observations at adjacent time points to create a set  $\mathscr{K}$  of design time points that are clinically interpretable, see Wu and Tian (2018, Sec. 12.2). We consider the age range T = [9.00, 19.00)and specify four design time points at each age, yielding K = 40 equally spaced age bins  $[9.00, 9.25), \ldots, [18.75, 19.00)$  corresponding to the design time points  $\mathscr{K} = \{9.00, 9.25, \ldots, 18.75\}$ . If the *i*th girl is observed within age bin  $[t_{(j)}, t_{(j+1)})$ , her corresponding design time point is  $t_{(j)}$ . This age binning has adequate clinical interpretation for pre-teens and adolescents, because two girls born within three months of each other have approximately the same age.

Let Y(t) be a girl's systolic blood pressure (SBP) at time point t years of age, and let  $A_1(t_1) = (y_{0.8}(t_1), \infty)$  and  $A_2(t_2) = (y_{0.8}(t_2), \infty)$  be the 80th percentile SBP ranges at ages  $t_1$  and  $t_2$  years, respectively. We estimate the RTPs  $RTP_{A_1, A_2}(t_2-3, t_2)$  and  $RTP_{A_1, A_2}(t_1, t_2)$  at a sequence of  $t_2$  values with  $t_1 = t_2-3$ and  $t_1 = 10$ , respectively, using both the copula-based smoothing method and the unstructured smoothing method for Caucasian and African American girls. These RTPs give quantitative measures for the probability of a girl's SBP being above the 80th percentile at age  $t_2$  years, given that her SBP is known to be above the 80th percentile at age  $t_1$  years. Further details about clinical interpretations of  $RTP_{,,.}(t_2-3, t_2)$  and  $RTP_{,,.}(t_1, t_2)$  for epidemiology studies can be found in Wu and Tian (2018, Chap. 12). We compute both the copula-based and the unstructured smoothing estimates and their corresponding bootstrap 95% empirical quantile point-wise confidence intervals using a resampling-subject bootstrap with 500 replications, following Wu and Tian (2018, Sec. 12.3.6).

Figure 3 shows the estimates of  $RTP_{0.8,0.8}(t_1, t_2)$  and their bootstrap 95% point-wise confidence intervals for Caucasian and African American girls over different age ranges. The top panels of Figure 3 show the estimated probability that, given that a Caucasian (Figure 3a) or an African American (Figure 3b) girl's SBP was higher than the age-specific 80th percentile at ages 9 to 14.5 years, her SBP is also higher than the age-specific 80th percentile three years later. The bottom panels of Figure 3 show the estimated probability that a Caucasian (Figure 3c) or an African American (Figure 3d) girl's SBP is higher than the age-specific 80th percentile three years later. The bottom panels of Figure 3 show the estimated probability that a Caucasian (Figure 3c) or an African American (Figure 3d) girl's SBP is higher



Figure 3. The estimated  $RTP_{\alpha_1,\alpha_2}(t_1,t_2)$  curves with  $\alpha_1 = \alpha_2 = 0.8$  using the copula-based smoothing estimator and the unstructured smoothing estimator, and the corresponding bootstrap 95% empirical quantile point-wise confidence intervals.

than the age-specific 80th percentile at ages 13 to 17.5 years, given that her SBP was already higher than the 80th percentile at age 10 years. All four panels of Figure 3 suggest that the RTPs for both Caucasian and African American girls vary between approximately 40% and 45% for different age ranges, which are much higher than the anticipated value of 20% if there were no tracking ability for the SBP of this population. These results, which are consistent with the findings reported in Wu and Tian (2013), indicate that the percentile SBP values for a girl at different ages have positive tracking ability, and hence are positively correlated.



Figure 4. The estimated  $RTP_{A_1, A_2}(10, t_2)$  curves with  $A_1(t) = A_2(t) = (115, \infty)$  using the copula-based smoothing estimator and the unstructured smoothing estimator, and the corresponding bootstrap 95% empirical quantile point-wise confidence intervals.

Figure 4 shows the estimated  $RTP_{A_1,A_2}(t_1,t_2)$  with  $A_1(t_1) = (115,\infty)$  and  $A_2(t_2) = (115,\infty)$  at a sequence of  $t_2$  values with  $t_1 = t_2 - 3$  and  $t_1 = 10$ , respectively, using the copula-based smoothing method and the unstructured smoothing method. Because  $A_1(t_1)$  and  $A_2(t_2)$  depend on the fixed SBP level of 115 mmHg, the estimated RTP curves in Figure 4a to Figure 4d are all increasing with age, suggesting that girls' SBP levels are increasing with age. These findings are consistent with those observed in Figure 3.

Comparing the bootstrap 95% point-wise confidence intervals in Figures 3 and 4, we observe that, for all the panels in Figure 3a to Figure 3d and Figure 4a to Figure 4d, the widths of the confidence intervals for the copula-based smoothing

estimates are narrower than those for the unstructured smoothing estimates. These results are similar to the simulation results summarized in Tables 1 and 2 and Figures 1 and 2, where we found that the copula-based smoothing estimator has smaller standard errors, in general, than those of the unstructured smoothing estimator.

## 6. Discussion

The RTP has been shown to be a useful measure of the tracking abilities of time-varying disease status and risk factors in longitudinal studies. We have proposed a copula-based smoothing method for estimating RTPs in longitudinal studies, either without covariates or with time-invariant categorical covariates. Our theoretical and simulation results demonstrate that the proposed method has major advantages over the current unstructured smoothing method. The proposed smoothing method consists of two steps: (a) computing the raw estimates at "design time points" based on a known copula model or a set of candidate copula models; (b) obtaining the functional RTP estimates at any time point by smoothing the raw estimates using a kernel smoothing method. We provide theoretical justifications for the proposed method by deriving the asymptotic mean squared errors of the estimators, and using a simulation study to demonstrate its finite-sample superiority over the unstructured smoothing method under the robust scenario that the copula model is not known, but is selected from a set of candidate copula models. Furthermore, our application to NGHS SBP data demonstrates that this copula-based smoothing method leads to clinically meaningful results in epidemiological studies.

A limitation of the proposed estimation method is that it applies only when there is no covariate. When including time-varying or continuous covariates, we require additional modeling assumptions for the copula models so that we can specify their dependence structures on the covariates. Further research is warranted to establish a flexible and clinically meaningful copula model that includes time-varying and continuous covariates. In addition, because the RTPs are often sensitive to the tail dependence structures, the copula models considered here may not be sufficient to handle all possible situations in practice. Therefore, it is preferable to extend our method to include more flexible models for the distribution functions. Examples of such an extension may include mixtures of copulas and vine copulas. Finally, the RTPs considered both here and by Wu and Tian (2013, 2018) are defined at two different time points. Given that there is practical interest in studying pediatric blood pressure distributions at more than three time points (NHBPEP (2004)), other statistical models for RTPs at three or more time points deserve future research.

## Supplementary Material

The online Supplementary Material provides detailed proofs for Theorems 1–4.

#### Acknowledgments

The NGHS data are available upon request from the BioLINCC website (https://biolincc.nhlbi.nih.gov/studies/nghs/) of the National Heart, Lung, and Blood Institute. The authors thank the investigators and participants of the NGHS study. The views expressed in this manuscript are those of the authors and do not necessarily represent the views of the NHLBI, the National Institutes of Health, or the US Department of Health and Human Services. The research of X. Y. Zhang and M. X. Wu was supported by the National Natural Science Foundation of China (No. 11771032 and 12271034). The authors are grateful to the editors and two referees for their helpful comments and suggestions.

#### References

- Chen, S. X. and Huang, T.-M. (2007). Nonparametric estimation of copula functions for dependence modelling. *Canadian Journal of Statistics* **35**, 265–282.
- Fan, J., Huang, T. and Li, R. (2007). Analysis of longitudinal data with semiparametric estimation of covariance function. *Journal of the American Statistical Association* 102, 632–641.
- 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.
- Hall, P., Racine, J. and Li, Q. (2004). Cross-validation and the estimation of conditional probability densities. *Journal of the American Statistical Association* 99, 1015–1026.
- Hall, P., Wolff, R. C. and Yao, Q. (1999). Methods for estimating a conditional distribution function. Journal of the American Statistical Association 94, 154–163.
- Härdle, W., Müller, M., Sperlich, S., and Werwatz, A. (2004). Nonparametric and Semiparametric Models. Springer Science & Business Media.
- Hoover, D. R., Rice, J. A., Wu, C. O. and Yang, L.-P. (1998). Nonparametric smoothing estimates of time-varying coefficient models with longitudinal data. *Biometrika* 85, 809– 822.
- Joe, H. (2014). Dependence Modeling with Copulas. Chapman and Hall/CRC.
- National Heart, Lung and Blood Institute Growth and Health Research Group (NGHSRG) (1992). Obesity and cardiovascular disease risk factors in black and white girls: The NHLBI growth and health study. *American Journal of Public Health* **82**, 1613–1620.
- National High Blood Pressure Education Program Working Group on High Blood Pressure in Children and Adolescents (NHBPEP Working Group) (2004). The fourth report on the diagnosis, evaluation, and treatment of high blood pressure in children and adolescents. *Pediatrics* 114, 555–576.
- Obarzanek, E., Wu, C. O., Cutler, J. A., Kavey, R.-E. W., Pearson, G. D. and Daniels S. R. (2010). Prevalence and incidence of hypertension in adolescent girls. *The Journal* of *Pediatrics* 157, 461–467.

- Rice, J. A. and Wu, C. O. (2001). Nonparametric mixed effects models for unequally sampled noisy curves. *Biometrics* 57, 253–259.
- Scaillet, O. and Fermanian, J.-D. (2002). Nonparametric estimation of copulas for time series. FAME Research Paper 57. Web: http://dx.doi.org/10.2139/ssrn.372142.
- Sklar, A. (1959). Fonctions de réprtition à n dimensions et leurs marges. Publications de Institut de Statistique de Université de Paris 8, 229–231.
- Wu, C. O. and Tian, X. (2013). Nonparametric estimation of conditional distribution functions and rank-tracking probabilities with longitudinal data. *Journal of Statistical Theory and Practice* 7, 259–284.
- Wu, C. O. and Tian, X. (2018). Nonparametric Models for Longitudinal Data: With Implementation in R. Chapman and Hall/CRC.

Xiaoyu Zhang

College of Statistics and Data Science, Beijing University of Technology, Beijing 100124, China. E-mail: zhangxiaoyu006@126.com

Mixia Wu

College of Statistics and Data Science, Beijing University of Technology, and Beijing Institute for Scientific and Engineer Computing, Beijing 100124, China.

E-mail: wumixia@bjut.edu.cn

Colin O. Wu

Office of Biostatistics Research, National Heart, Lung and Blood Institute, Bethesda, MD 20892, USA.

E-mail: wuc@nhlbi.nih.gov

(Received December 2021; accepted September 2022)