

## Kernel Ridge Estimator for the Partially Linear Model under Right-Censored Data

Syed Ejaz Ahmed <sup>1</sup>, Dursun Aydın <sup>2</sup>, and Ersin Yılmaz <sup>2</sup>

<sup>1</sup> Brock University, Faculty of Mathematics and Science, Department of Mathematics and Statistics, Niagara Region, 1812 Sir Isaac Brock Way, St. Catharines, ON, L2S 3A1, Canada.

<sup>2</sup> Mugla Sitki Kocman University, Faculty of Science, Department of Statistics, 48000, Mugla, Turkey.

Received: 08/12/2020, Revision received: 29/01/2021, Published online: 03/04/2021

**Abstract.** *Objective:* This paper aims to introduce a modified kernel-type ridge estimator for partially linear models under randomly-right censored data. Such models include two main issues that need to be solved: multi-collinearity and censorship. To address these issues, we improved the kernel estimator based on synthetic data transformation and kNN imputation techniques. The key idea of this paper is to obtain a satisfactory estimate of the partially linear model with multi-collinear and right-censored using a modified ridge estimator. *Results:* To determine the performance of the method, a detailed simulation study is carried out and a kernel-type ridge estimator for PLM is investigated for two censorship solution techniques. The results are compared and presented with tables and figures. Necessary derivations for the modified semiparametric estimator are given in appendices.

**Keywords.** Kernel Smoothing, KNN Imputation, Multi-Collinear Data, Partially Linear Model, Ridge Type Estimator, Right-Censored Data.

---

Corresponding Author: Syed Ejaz Ahmed (sahmed5@brocku.ca)  
Dursun Aydın (duaydin@hotmail.com)  
Ersin Yılmaz (yilmazersin13@hotmail.com)

MSC: 62G05, 62G08, 62N01.

## 1 Introduction

A partially linear model includes multicollinearity as follows:

$$y_i = \sum_j^p x_{ij}\beta_j + g(z_i) + \varepsilon_i, \quad 1 \leq i \leq n, \quad (1.1)$$

where  $y_i$  represents the response values,  $x_{ij}$  denotes the values of covariates correlated with each other,  $\beta_j$  represents unknown regression coefficients for the parametric component,  $g(z_i)$  is the nonparametric part of the model with single nonparametric covariate  $z_i$  and an unknown smooth function  $g(\cdot)$ , and  $\varepsilon_i$  denotes the random error terms with constant variance  $\sigma_\varepsilon^2$  and mean zero. To simplify this notation, the matrix and vector form of the model (1.1) can be written as follows

$$\mathbf{y} = \mathbf{X}^T \boldsymbol{\beta} + \mathbf{g} + \boldsymbol{\varepsilon}, \quad (1.2)$$

where  $\mathbf{y} = (y_1, y_2, \dots, y_n)^T$  is the  $(n \times 1)$  vector of response values,  $\mathbf{X}^T = (x_1, \dots, x_p)$  is the  $(n \times p)$  covariate matrix formed by observations of independent variables and assumed to be a full-ranked matrix,  $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)^T$  is the  $(p \times 1)$  vector of regression coefficients to be estimated,  $\mathbf{g} = (g(z_1), g(z_2), \dots, g(z_n))^T$  is the  $(n \times 1)$  dimensional vector of nonparametric component, and  $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$  is the  $(n \times 1)$  vector of random error terms with  $E(\boldsymbol{\varepsilon}) = 0$  and  $E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^T) = \sigma_\varepsilon^2 \mathbf{I}_n$ . A number of researchers have studied the fitness of the model defined in (1.2). Examples of such studies include Schimek (2000), Ruppert et al. (2003), and Liang (2006).

The paper discusses two important problems that arises in model (1.2), multicollinearity and censored response values.

Multi-collinearity is an increasingly common problem in recent years, as data collection methods have been developed and data structures have evolved into high dimensionality. It is well-established that standard regression phenomenon assumes that covariates are linearly uncorrelated. However, in many application areas such as biology, economics, industry, finance, medical studies, and especially recent work in bioinformatics, this assumption often breaks down because there are many more explanatory variables to model the response variable. Therefore, multi-collinearity

is an inevitable problem for any linear, nonparametric or semi-parametric modeling procedures.

It is quite common to have observations with missing values for response variables in applied fields such as survival analysis. These missing observations are sometimes unavoidable, due to the fact that individuals can withdraw from a study at any time. In general, censoring occurs when there is incomplete information about the survival time of some individuals in the study. The censored data problem, where the observed value of a variable is partially known, is related to the missing data problem. Ordinary regression techniques cannot be used directly in the fitting procedure of model (1.1) because censored data leads to biased estimates. In such cases, one typical approach is to impute, or "fill in", the missing observations; another is to transform response observations.

It should be emphasized that when there are many independent variables affecting the right-censored response observations, the multi-collinearity problem may arise in the setting of regression analysis. Specifically, this article seeks a solution for situations where both right-censorship and collinearity appear simultaneously in the data set. This is because, based on our experience, datasets have frequently contained both these problems in recent years.

This paper solves the multi-collinearity problem by using the kernel-type ridge estimator for partially linear models studied by Yüzbaşı et al. (2017). To overcome censorship, we use two methods: the synthetic data transformation method described in Koul et al. (1981) and the kNN imputation method studied by Batista and Monard (2002) and Ahmed et al. (2020). Note also that kernel ridge regression under censored data has been previously studied by Shim (2005) in a nonparametric setting. As is known, in ridge-type estimators, the shrinkage parameter plays a crucial role in accuracy of the estimation. Moreover, the bandwidth parameter of the kernel smoothing method controls the amount of penalty term in minimization criterion given in equation (3.1). Determining both the ridge and bandwidth parameters are therefore two additional issues covered in this paper. To achieve improved  $AIC_c$ , the criterion proposed by Hurvich et al. (1998) is used in the determination of both parameters.

In the context of right-censored data, several authors have studied estimation of the partially linear model (1.2). Some of these include Orbe et. al. (2003) and Liang and Zhou (2008). In addition, the estimation of right-censored response values was proposed by Kaplan-Meier (1958), and later Miller (1976) proposed Kaplan-Meier (K-

M) weights for the linear regression model estimation using a K-M estimator.

According to given information above, the contributions of this paper can be summarized as follows:

- A modified kernel-type ridge estimator is introduced to estimate the partially linear model with right-censored data.
- Two solution techniques with different fundamentals are used to solve the censorship problem. One of these is a synthetic data transformation based on the Kaplan-Meier estimator of distribution of censoring variable (Koul et al., 1981). The other is kNN imputation, a machine learning method. It is a fully nonparametric model that provides different perspectives on the estimation process.

The paper is arranged as follows: In section 2, the right-censored data concept, synthetic data transformation and kNN imputation methods are introduced with their important aspects. Section 3 presents the ridge-type kernel estimator for a censored, partially linear model based on both synthetic data and kNN imputation. In Section 4, some performance evaluation metrics are defined. A detailed Monte-Carlo simulation study and results are given in Section 5. Finally, concluding remarks are made in Section 6.

## 2 Right-Censored Data and Solution Techniques

This paper focuses on estimating the vector  $\beta$  of regression coefficients and the unknown smooth function  $g$  in model (1.2), while the response values of  $y_i$  are observed incompletely and censored from the right by a random censoring variable  $c_i$ ,  $i = 1, \dots, n$ . However,  $x_{ij}$  and  $z_i$  variables are completely observed. Consequently, right-censored triplets  $\{x_{ij}, z_i, y_i\}$  are turned into incomplete observations  $\{x_i, z_i, t_i, \delta_i\}_{(i=1)}^n$  as follows:

$$t_i = \min(y_i, c_i) \text{ and } \delta_i = I(y_i \leq c_i), 1 \leq i \leq n, \quad (2.1)$$

where  $t_i$  represents the right-censored response values that are updated according to existence of the censorship with distribution  $J$ ,  $c_i$  denotes the values of the censoring variable with distribution  $G$  and  $\delta_i$  values are binary scores of the censoring indicator, which carries the information of the existence of the censorship. Note also that it is assumed that  $y_i$  and  $c_i$  are independent random variables with unknown distributions  $F$  and  $G$ , respectively.

Now, it can be clearly seen that model (1.2) cannot be estimated directly using an ordinary semiparametric modelling procedure due to incomplete responses. Although

$t_i$ 's are obtained according to censorship, they do not involve the effect of censored observations which cause biased estimates. In order to overcome this problem, there are some solution methods in the literature. The most important methods can be divided into three categories: transformation (Koul et al. 1981), imputation (Batista and Monard 2002; Yenduri and Iyengar, (2007), and weights (Miller, 1976; Stute, 1993; Orbe et al. 2003). In this paper, the first two types are considered because, for the purposes of this paper, synthetic data transformation and imputation techniques are far easier, in terms of computation and formulation, than a weights-based method.

In the analysis of right-censored data, there are two fundamental assumptions (A1 and A2) which provide consistent and accurate estimates. These assumptions have previously been described in detail by Stute (1993), Koul et al. (1981) and Miller (1976). Therefore, they will be introduced only cursorily here:

**A1.** Completely observed  $y_i$ 's and  $c_i$ 's should be independent.

**A2.**  $P(y_i < c_i | y_i, x_i, z_i) = P(y_i \leq c_i | y_i)$ .

A1 is the ordinary assumption for survival analysis and allows acquisition of a meaningful model. If A1 is broken, then to obtain an accurate model, significantly more information is needed. On the other hand, A2 means that given time of failure, values of explanatory variables cannot provide any further information about response data points. Note that both A1 and A2 cannot assume the independency between explanatory variables, but in classical regression phenomenon, correlation between  $x_i$ 's is not preferred, as explained in the introduction section. It should be emphasized that this paper's interests lie in two problems: right-censored responses and multi-collinearity. This section prepares us to handle censorship problem. Multi-collinearity under censorship is discussed in Section 3.

## 2.1 Synthetic Data Transformation

Synthetic data transformation is a widely used technique to overcome censorship problem in the data in the modelling context. There are a number of modified synthetic data transformations in the literature, some of which are mentioned in Section 1. However, this paper will obtain synthetic responses based on Koul et al. (1981). Accordingly, data transformation is computed by

$$t_{iG} = \frac{\delta_i t_i}{1 - G(t_i)} = \frac{\delta_i t_i}{\bar{G}(t_i)}, \quad (2.2)$$

where  $G$  is the distribution of  $c_i$ , as explained before. After computation of the synthetic responses, model (1.2) can be rewritten with synthetic data as follows:

$$t_G = \mathbf{X}^T \boldsymbol{\beta} + \mathbf{g} + \boldsymbol{\varepsilon}_G, \quad (2.3)$$

where  $\boldsymbol{\varepsilon}_G = (\varepsilon_{1G}, \dots, \varepsilon_{nG})$  is a vector of random error terms for known censoring distribution  $G$  and  $E(\boldsymbol{\varepsilon}_G) = 0$  (see, for example, Aydın and Yılmaz, 2018, for a more detailed discussion). Because  $G$  is generally unknown and has to be estimated,  $G$  is replaced by its Kaplan-Meier estimation, as introduced by Koul et al. (1981), below:

$$\hat{G}(s) = 1 - \prod_{i=1}^n \left( \frac{n-i}{n-i+1} \right)^{I_{[t_{(i)} \leq s, \delta_{(i)}=0]}}, \quad (s \geq 0), \quad (2.4)$$

where  $\left\{ (t_{(i)}, \delta_{(i)}) \right\}_{i=1}^n$  denotes ordered data pairs associated with  $t_{(i)}$ . Note that Koul et al. (1981) shows that for  $n \rightarrow \infty$ ,  $\hat{G}(s) \rightarrow G(s)$ . Moreover, the main idea of using synthetic values is that  $t_{iG}$ 's have exactly the same expected value as the observed response variable:  $E(t_{i\hat{G}}) \cong E(t_{iG}) \cong E(y_i)$ , as discussed by Aydın and Yılmaz (2018).

## 2.2 The kNN Imputation Technique

This section has been prepared to introduce the kNN imputation technique to impute estimated scores instead of censored data points independent of distribution, which is the most important difference between the kNN imputation from synthetic data transformation. Advantages and disadvantages of this technique are given as follows:

---

### Advantages

- Imputed values are derived from actual values not synthetic or constructed values
- kNN provides additional information using explanatory variables.
- kNN imputation is a fully nonparametric method and it does not make any assumption about the relationship between  $x$  and  $y$  or  $x$  and  $z$ .
- kNN can work for both discrete and continuous attributes. For discrete attributes, it uses the most frequently used value among  $k$ -nearest neighbors.
- For continuous attributes, it uses mean value of  $k$ -nearest neighbors

---

### Disadvantages

- It cannot be guaranteed to obtain the exact same expected values of imputed dataset ( $y_i^k$ ) and completely observed response  $y_i$
  - It is possible that the technique does not produce consistent results. Thus, estimates do not always get better when the sample size is larger.
  - Statistical properties such as bias and variance of the obtained estimator cannot be computed, because it is a fully nonparametric method.
-

The kNN is a similarity-based technique, that is, it depends on the distance between observations. There are various metrics to evaluate these distances. However, in this paper, the Euclidean norm, which is widely used in the literature, is preferred. Accordingly, computation of the Euclidean norm is given by

$$d(a, b) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}, \quad (2.5)$$

where  $d(a, b)$  represents the value of the distance. An algorithm is developed to perform the kNN imputation for right-censored observations and is given in Algorithm 1. This algorithm is also used in Ahmed et al. (2020).

---

**Algorithm 1** Algorithm of kNN imputation

---

- 1: **Input:** Right-censored dataset  $t_i$
  - 2: Censorship indicator  $\delta_i$
  - 3: Number fo nearest neighbours  $k$
  - 4: Values of explanatory variable  $z_i$
  - 5: **Output:** Imputed dataset  $\mathbf{y}^{kmn} = (y_1^k, \dots, y_n^k)^\tau$
  - 6: **begin**
  - 7: **for** ( $i = 1$  to  $n$ )
  - 8: **if** ( $\delta_i = 0$ ) **do** Step 4 (if data point is censored)
  - 9: **for** ( $j = 1$  to  $n$ )
  - 10: Find the Euclidean distances (2.5) between  $z_j$  and  $z_i$  for each censored data point
  - 11: **end** (If statement in Step 3)
  - 12: sort the distances from small to large
  - 13: **end** (for loop in Step 2)
  - 14: **end** (for loop in Step 4)
  - 15: **for** ( $j = 1$  to  $k$ )
  - 16: Take the first *uncensored*  $k$  values of  $t_i$  associated to sorted distances
  - 17: Calculate the  $i^{\text{th}}$  imputed value ( $y_i^k$ ) with average of nearest  $k$  records of  $y_i$
  - 18: Replace the imputed values ( $y_i^k$ ) with censored data points ( $z_i, \delta_i = 0$ ) in censored data set  $\mathbf{t} = (t_1, \dots, t_n)^T$
  - 19: **end** (for loop in Step 10)
  - 20: Return  $\mathbf{y}^k = (y_1^k, \dots, y_n^k)^T$
  - 21: **end.**
- 

See Ahmed et al. (2020) for a more detailed discussion on the kNN imputation

method based on right-censored data.

### 3 Censored Ridge-Type Kernel Estimator (CRK)

We first consider the nonparametric estimation of the unknown function  $g(z)$  in (1.1) based on uncensored data. For simplicity, we assume that  $\beta$  in model (1.1) is known. Accordingly, the relationship between  $y_i - \sum_j^p x_{ij}\beta_j$  and  $z_i$  can be defined as

$$y_i - \sum_j^p x_{ij}\beta_j = g(z_i) + \varepsilon_i, i = 1, 2, \dots, n, \quad (3.1)$$

The purpose here is to approximate  $g(z_i)$  closely. In this sense, an estimator is discussed by Nadaraya and Watson (1964), and it is also referred to as the kernel estimator, given by

$$g(z) = \sum_{i=1}^n w_{ih}(z_i) (y_i - x_i\beta)^2 = \mathbf{W}_h (\mathbf{y} - \mathbf{X}^T \beta) = \hat{g}, \quad (3.2)$$

where  $h$  is a bandwidth (or smoothing) parameter,  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  are known  $p$ -vectors, as defined in (1.1), and  $\mathbf{W}_h$  is a kernel smoothing matrix with  $i^{\text{th}}$  entries  $w_{hi}$ , defined by

$$w_{hi} = K\left(\frac{z - z_i}{h}\right) / \sum_{i=1}^n K\left(\frac{z - z_i}{h}\right) = \frac{K(u)}{\sum K(u)} = \mathbf{W}_h, \quad (3.3)$$

where  $K(u)$  is a kernel function satisfying these properties: (i)  $\int K(u)du = 1$ , (ii)  $K(u) = K(-u)$ . The main characteristic of the kernel function is that it gives more weight to observations near  $z$  and less weight to observations far from  $z$ . Also,  $K(u)$  controls the shape of the estimated nonparametric curve with bandwidth parameter  $h$ . For this study, bandwidth is determined by GCV criterion.

Using the equation (3.2) and the matrix-vector form of model (3.1), we can get the following partial residual:

$$\varepsilon = \mathbf{y} - \mathbf{X}^T \beta - \hat{g} = (\mathbf{I}_n - \mathbf{W}_h) (\mathbf{y} - \mathbf{X}^T \beta) = (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \beta), \quad (3.4)$$

where  $\tilde{\mathbf{X}} = (\mathbf{I}_n - \mathbf{W}_h) \mathbf{X}$  and  $\tilde{\mathbf{y}} = (\mathbf{I}_n - \mathbf{W}_h) \mathbf{y}$ . We thus obtain a transformed data set based on kernel residuals. Considering the partial kernel residuals given in (3.4), the estimator of the parameters vector  $\beta$  can be obtained by minimizing the weighted

residual sum of squares:

$$\begin{aligned} L_s(\boldsymbol{\beta}, h) &= \left[ (\mathbf{I}_n - \mathbf{W}_h)(\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta}) \right]^T \left[ (\mathbf{I}_n - \mathbf{W}_h)(\mathbf{y} - \mathbf{X}^T \boldsymbol{\beta}) \right] \\ &= \left( \tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \boldsymbol{\beta} \right)^T \left( \tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \boldsymbol{\beta} \right). \end{aligned} \quad (3.5)$$

Differentiating with respect to  $\boldsymbol{\beta}$  of (3.1), we obtain the normal equations by

$$\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\beta} = \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}. \quad (3.6)$$

To solve the normal equations for  $\boldsymbol{\beta}$ , we multiply both sides of (3.3) by  $\left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1}$ . Thus, the weighted least square estimator of  $\boldsymbol{\beta}$  is defined as

$$\widehat{\boldsymbol{\beta}} = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}. \quad (3.7)$$

Using equation (3.2) for the unknown regression function vector  $\mathbf{g}$ , an updated equation can be written as follows:

$$\hat{\mathbf{g}} = \mathbf{W}_h \left( \mathbf{y} - \mathbf{X}^T \widehat{\boldsymbol{\beta}} \right). \quad (3.8)$$

It should be emphasized that equations (3.7) and (3.8) are considered as if no perfect or exact relationship between the columns of the matrix  $\mathbf{X}$  or the inverse matrix  $\left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1}$  exists. As noted in introduction, this paper is designed to overcome these problems. To this end, the standard remedy is to use a biased estimation method such as ridge regression, as proposed by Hoerl and Kennard (1970). For linear model  $\mathbf{y} = \mathbf{X}^T \boldsymbol{\beta} + \boldsymbol{\varepsilon}$ . For  $k > 0$ , the ridge estimate of the parameters vector  $\boldsymbol{\beta}$  in a linear regression model is given by

$$\widehat{\boldsymbol{\beta}}_R = \left( \mathbf{X}^T \mathbf{X} + k \mathbf{I}_p \right)^{-1} \mathbf{X}^T \mathbf{y}, \quad (3.9)$$

where  $\mathbf{I}_p$  is  $(p \times p)$  dimensional identity matrix and  $k$  is a ridge (or shrinkage) parameter to be selected using plug-in methods or common criteria such as generalized cross validation (GCV), Akaike information criterion, and so on. In this paper, it is selected by GCV criterion.

To fit model (1.1) to the data, a ridge procedure can be associated with the idea of hints due to Speckman (1988), where the ridge estimator is the minimizer of the

penalized least squares criterion

$$\begin{aligned} L_R(\boldsymbol{\beta}, \mathbf{g}; k) &= \sum_{i=1}^n \left( \tilde{y}_i - \mathbf{X}_i^T \boldsymbol{\beta} - \mathbf{g}(z_i) \right)^2 + k \sum_{j=1}^p \beta_j^2 \\ &= \sum_{i=1}^n \left( \tilde{y}_i - \mathbf{X}_i^T \boldsymbol{\beta} - \mathbf{g}(z_i) \right)^2 + \sum_{j=1}^p (0 - k\beta_j)^2. \end{aligned} \quad (3.10)$$

An equivalent way to write the criterion in (3.10) in matrix form is

$$L_R(\boldsymbol{\beta}, \mathbf{g}; k) = \left( \tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \boldsymbol{\beta} - \mathbf{g} \right)^T \left( \tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \boldsymbol{\beta} - \mathbf{g} \right) + k \|\mathbf{0} - \boldsymbol{\beta}\|^2. \quad (3.11)$$

From (3.11), the ridge-type kernel estimator of  $\boldsymbol{\beta}$  and  $\mathbf{g}$  are defined respectively by

$$\widehat{\boldsymbol{\beta}}_R = \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k\mathbf{I}_p \right)^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}, \quad (3.12)$$

and

$$\widehat{\mathbf{g}}_R = \mathbf{W}_h \left( \mathbf{y} - \mathbf{X}^T \widehat{\boldsymbol{\beta}}_R \right). \quad (3.13)$$

See Aydın and Yılmaz (2018) for proofs, the statistical properties of the estimators given in (3.11), (3.12) and other details about using a ridge-type kernel estimator with a partially linear model.

In this section, a ridge-type kernel estimator is introduced to estimate right-censored responses under collinear data. As stated in Section 2, the problem of censored response observations is solved using two different techniques. We therefore have two new response variables,  $t_{i\hat{G}}$  and  $y_i^k$ , obtained by synthetic data and imputation techniques, respectively. The core idea presented here is to replace the response vector  $\mathbf{y}$  with the vector  $t_{i\hat{G}}$  from synthetic data and the vector  $\mathbf{y}^k$  obtained by kNN imputation method. Thus two different ridge-type kernel estimators will be obtained for partially linear model (1.2). With this context in mind, the following theorems establishes these procedures:

**Theorem 3.1.** (CRK with synthetic data  $t_{i\hat{G}}$ ): Let  $\mathbf{t}_{\hat{G}} = \tilde{\mathbf{X}}^T \boldsymbol{\beta} + \tilde{\boldsymbol{\varepsilon}}$  where  $\tilde{\boldsymbol{\varepsilon}} = \tilde{\mathbf{g}} + \tilde{\boldsymbol{\varepsilon}}_{\hat{G}}$ ,  $\tilde{\mathbf{g}} = (\mathbf{I}_n - \mathbf{W}_h) \mathbf{g}$  and  $\tilde{\boldsymbol{\varepsilon}}_{\hat{G}} = (\mathbf{I}_n - \mathbf{W}_h) \boldsymbol{\varepsilon}_{\hat{G}}$ . Also,  $\tilde{\mathbf{X}}$  is a  $n \times p$  matrix and  $\tilde{\mathbf{t}}_{\hat{G}} = (\mathbf{I}_n - \mathbf{W}_h) \mathbf{t}_{\hat{G}}$  in which  $\mathbf{t}_{\hat{G}}$  is a  $n \times 1$  synthetic data vector replaced with  $\mathbf{y}$ . If  $\mathbf{W}_h$  is an arbitrary smoother matrix then the ridge regression estimates may be calculated by augmented data

$$\mathbf{X} = \begin{pmatrix} \tilde{\mathbf{X}} \\ \sqrt{k}\mathbf{I}_p \end{pmatrix} \text{ and } \mathbf{t}_{\hat{G}} = \begin{pmatrix} \tilde{\mathbf{t}}_{\hat{G}} \\ \mathbf{0}_p \end{pmatrix}. \quad (3.14)$$

Then kernel-type ridge estimator for  $\beta$  is denoted as  $\hat{\beta}_{\hat{G}}$  and given by

$$\hat{\beta}_{\hat{G}} = (\tilde{X}^T \tilde{X} + kI_p)^{-1} \tilde{X}^T t_{\hat{G}}. \quad (3.15)$$

Proof of the Theorem 3.1 is given in Appendix A1.

It should be noted that because of  $t_{i\hat{G}} \rightarrow t_{iG}$  and  $E(t_{iG}) \cong E(y_i)$  when  $n \rightarrow \infty$ . Thus, assumption of  $E(\varepsilon_{i\hat{G}}) \cong E(\varepsilon) = 0$  will be ensured. Note also that when  $k = 0$ , the ridge-type estimate reduces to the estimate problem in the equation (3.7) (see Speckman, 1988). Also, it is seen that there is a formal similarity between the equation (3.11) and ridge estimator of the linear regression model. Combining equations (3.11) and (3.13) we get the estimator of the vector  $g$  as

$$\hat{g}_{\hat{G}} = W_h (\tilde{t}_{\hat{G}} - X^T \hat{\beta}_{\hat{G}}). \quad (3.16)$$

**Theorem 3.2.** (CRK with imputed  $y^k$ ): Similar to Theorem 3.1, if  $y_i$  is replaced by imputed values  $y_i^k$ 's using kNN imputation, then ridge-type estimator for  $\beta$  and  $g$  in a partially linear model can be defined, respectively, as follows:

$$\hat{\beta}_k = (\tilde{X}^T \tilde{X} + kI_p)^{-1} \tilde{X}^T \tilde{y}^k, \quad (3.17)$$

and

$$\hat{g}_k = W_h (\tilde{y}^k - X^T \hat{\beta}_k), \quad (3.18)$$

where  $\tilde{X} = (I_n - W_h) X$ , and  $\tilde{y}^k = (I_n - W_h) y^k$ . Proof of equations (3.17)-(3.18) are given in Appendix A2.

Theoretically,  $E(y^k) \cong E(y)$  for  $n \rightarrow \infty$  cannot be proven for Theorem 3.2, because, kNN has a fully nonparametric nature. However, in practice, it can be seen that the larger the sample size, the better the results, although this improvement in the estimates cannot be guaranteed. For our results, see Section 5.

## 4 Performance Measures

This section prepared to introduce evaluation metrics of the proposed modified estimators. Bias and variance of the regression coefficients  $\beta_j$ ,  $j = 1, 2, \dots, p$ , variance of errors ( $\hat{\sigma}_{\varepsilon}^2$ ), mean square error (MSE), and relative efficiencies (RE) are considered to determine the superiority of the estimators from the parametric and nonparametric components of the partially linear model. Note that in calculation of the metrics, for simplicity, we use only notation  $y_i$  to present response variable. One can use  $t_{i\hat{G}}$  or  $y_i^k$ , instead of  $y_i$ , for synthetic data and kNN imputation methods, respectively.

#### 4.1 Measures for Parametric Component

Equations (3.14)-(3.17) are rewritten as follows for  $y_{iG}$  equivalently:

$$\widehat{\beta}_r = \left[ I_p + k (\tilde{X}^T \tilde{X})^{-1} \right]^{-1} \beta. \quad (4.1)$$

It can be seen from (4.1) that the ridge-type estimator  $\widehat{\beta}_r$  is an obviously biased estimator because of  $\left[ I_p + k (\tilde{X}^T \tilde{X})^{-1} \right]^{-1} \neq I_p$ . Note that shrinkage parameter  $k \geq 0$  controls the size of estimated coefficients  $\widehat{\beta}_r$  and, consequently, the bias of the estimates (See Hoerl and Kennard, 1970). From that, we can surmise that if  $R_r = \left[ \tilde{X}^T \tilde{X} + kI_p \right]^{-1}$ , then the expected value, mean, and variance of  $\widehat{\beta}_r$  can be obtained as follows:

$$E(\widehat{\beta}_r) = R_r (\tilde{X}^T \tilde{X} \beta + \tilde{X}^T \tilde{g}) = \beta - kR_r \beta + R_r \tilde{X}^T \tilde{g}, \quad (4.2)$$

$$\text{Bias}(\widehat{\beta}_r) = B(\widehat{\beta}_r) = R_r \tilde{X}^T \tilde{g} - kR_r \beta, \quad (4.3)$$

$$\text{Variance}(\widehat{\beta}_r) = V(\widehat{\beta}_r) = \sigma^2 R_r \tilde{X}^T (I_n - W_h)^T (I_n - W_h) \tilde{X} R_r. \quad (4.4)$$

Note that the given equations are obtained using joint notation  $\widehat{\beta}_r$ . To compute moments (4.2)-(4.4) for  $\widehat{\beta}_{\hat{G}}$  in (3.14) or  $\widehat{\beta}_k$  in (3.17),  $\widehat{\beta}_r$  should be replaced by those notations.

#### 4.2 Measures for Nonparametric Component

The performance of the nonparametric component of the partially linear model is evaluated here by using mean squared error (MSE), which is a widely used metric in the nonparametric literature. MSE can be obtained for the synthetic data and the kNN imputation methods as follows:

$$MSE_{\hat{G}}(\hat{g}_{\hat{G}}, g) = E \left[ \left\{ \hat{g}_{\hat{G}}(z_i) - g(z_i) \right\}^2 \right] = \frac{1}{n} \sum_{i=1}^n \left[ g_{\hat{G}}(z_i) - g(z_i) \right]^2, \quad (4.5)$$

and

$$MSE_k(\hat{g}_k, g) = E \left[ \left\{ \hat{g}_k(z_i) - g(z_i) \right\}^2 \right] = \frac{1}{n} \sum_{i=1}^n \left[ \hat{g}_k(z_i) - g(z_i) \right]^2, \quad (4.6)$$

where  $MSE_{\hat{G}}$  and  $MSE_k$  denote the synthetic data and kNN methods, respectively.

### 4.3 Measures for the Model

In order to show the performance of the all estimated models using the mentioned methods, error variance of the models should be calculated. In this context, let  $\sigma_\varepsilon^2$  and  $\mathbf{y}$  be the joint notations. Accordingly, to estimate  $\sigma_\varepsilon^2$ , the residual sum of squares (RSS) are generally used. For the partially linear models, this can be obtained as follows:

$$\begin{aligned} \text{RSS}(\hat{\mathbf{y}}) &= (\hat{\mathbf{y}} - \mathbf{y})^T (\hat{\mathbf{y}} - \mathbf{y}) \text{ where } (\hat{\mathbf{y}} = \mathbf{X}^T \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\xi}} = \mathbf{H}_h \mathbf{y}) \\ &= (\mathbf{y} - \mathbf{H}_h \mathbf{y})^T (\mathbf{y} - \mathbf{H}_h \mathbf{y}) = \mathbf{y}^T (\mathbf{I}_n - \mathbf{H}_h) \mathbf{y}, \end{aligned} \quad (4.7)$$

where  $\mathbf{H}_h$  is known as a hat or smoother matrix, depending on bandwidth parameter  $h \geq 0$ , which is idempotent and used in the acquisition of fitted values for the linear models. Computation of  $\mathbf{H}_h$  is given by

$$\mathbf{H}_h = \mathbf{W}_h + (\mathbf{I}_p - \mathbf{W}_h) \tilde{\mathbf{X}} \mathbf{R}_r \tilde{\mathbf{X}}^T. \quad (4.8)$$

To obtain an estimation of  $\sigma_\varepsilon^2$  from that, the expected values of (4.7) are needed. This can be written as

$$E[\text{RSS}(\hat{\mathbf{y}})] = \sigma_\varepsilon^2 \left[ n - \text{tr} (2\mathbf{H}_h - \mathbf{H}_h^T \mathbf{H}_h) \right] + E(\mathbf{y}^T) (\mathbf{I}_n - \mathbf{H}_h)^2 E(\mathbf{y}), \quad (4.9)$$

where  $\sigma_\varepsilon^2 \left[ n - \text{tr} (2\mathbf{H}_h - \mathbf{H}_h^T \mathbf{H}_h) \right]$  denotes the variance and  $E(\mathbf{y}^T) (\mathbf{I}_n - \mathbf{H}_h)^2 E(\mathbf{y})$  represents the squared-bias of the model. In many applications,  $\sigma_\varepsilon^2$  is unknown. Because its estimate  $\hat{\sigma}_\varepsilon^2$  is used instead of  $\sigma_\varepsilon^2$ , it is calculated similar to OLS regression:

$$\hat{\sigma}_\varepsilon^2 = \frac{E(\text{RSS})}{\text{tr}[(\mathbf{I}_n - \mathbf{H}_h)^2]} = \frac{1}{n} \frac{\|(\mathbf{I}_n - \mathbf{H}_h) \mathbf{y}\|^2}{n - p}, \quad (4.10)$$

where  $\text{tr}[(\mathbf{I}_n - \mathbf{H}_h)^2] = (n - p)$  and  $p$  is the number of parameters. Note that  $(n - p)$  is the degree of freedom and it therefore should be obvious that (4.10) fits the classical variance definition. In this context, it can be said that  $\hat{\sigma}_\varepsilon^2$  is the bias estimate of  $\sigma_\varepsilon^2$ . This is a very common metric to evaluate the quality of estimated regression models.

As mentioned before, this paper considers two different estimators and equation (4.10) is therefore rearranged according to the kNN imputation and synthetic data methods with respect their response variables ( $\mathbf{y}^k$  and  $t_{\hat{\xi}}$ ) and corresponding hat matrices. Also note that, in order to compare the two mentioned methods in the ridge-type estimators under censored data context, a relative measure is needed. Therefore, associated with equation (4.10), a new measurement is introduced in Definition 4.1.

**Definition 4.1.** The relative efficiency based on the model variances  $\hat{\sigma}_{\varepsilon_G}^2$  and  $\hat{\sigma}_{\varepsilon_k}^2$  is specified as follows:

$$RE(t_{\hat{G}}, \hat{y}^k) = \frac{\left\| (I_n - H_{h_G}) t_{\hat{G}} \right\|^2}{\left\| (I_n - H_{h_k}) y^k \right\|^2} = \frac{\hat{\sigma}_{\varepsilon_G}^2}{\hat{\sigma}_{\varepsilon_k}^2}. \quad (4.11)$$

Note that if  $RE(t_{\hat{G}}, \hat{y}^k) < 1$ , then it can be said that fitted model based on the synthetic data technique is better than that of the kNN imputation.

## 5 Simulation Study

This section prepared to see behaviors of the introduced two CRK estimators under different conditions and scenarios. A detailed simulation experiments has been designed accordingly. Note that this paper has two main purposes: (i) to estimate right-censored data by partially linear models when multi-collinearity exists by using two different approaches, and (ii) to decide which approach is better under certain conditions and making inferences.

In order to obtain simulated data sets, a partially linear model is used and the values of the response variable are obtained using the following model:

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \beta_4 x_{4i} + g(z_i) + \varepsilon_i, \quad i = 1, 2, \dots, n, \quad (5.1)$$

where  $y_i$  represents completely observed values,  $(x_1, x_2, x_3, x_4)$  are correlated explanatory variables,  $\beta = (\beta_1, \beta_2, \beta_3, \beta_4)^T = (4, 2, 0.5, -1)^T$  and  $\varepsilon \sim N(0, \sigma^2 I)$ . Here, the simulation experiment is repeated 1000 times for  $\sigma^2 = 0.3$  and  $\sigma^2 = 1$  in three sample sizes:  $n = 30, 100$  and  $250$ . In addition, to add censorship to the response variables, the censoring variable is generated randomly by using distribution of  $y_i$ 's as  $c_i \sim N(\mu_y, \sigma_y^2)$ .  $I(c_i > y_i)$  where  $\mu_y$  and  $\sigma_y^2$  are the mean and variance of completely observed response variable  $y$ . Using equation (2.1), censored responses  $t_i$ 's are obtained. Note that censoring levels are decided as CL=(5%, 35%). The nonparametric component of the model is generated as follows:

$$g(z_i) = -0.1812 - 0.3221z_i + 4 \sin(z_i^2) + \exp(z_i), \quad (5.2)$$

where  $z_i = -2.4(i - 0.5)/n$  uniformly produced on interval  $[0, 1]$ . The covariates are obtained from normal distribution with two correlation levels. In this study, the level of multicollinearity ( $\rho$ ) is set as  $\rho = 0.85$  and  $\rho = 0.99$ . Note that we check existence of

the multicollinearity using a condition index ( $CI$ ), which provides the mightiness of  $\rho$  between variables.  $CI$  is widely used to measure collinearity (see Belsley et al., 1980). The formulation of  $CI$  is given by

$$CI = \sqrt{[\lambda_{\max}(X^T X) / \lambda_{\min}(X^T X)]}. \tag{5.3}$$

In addition, to evaluate the performances of the models and their components, measurement tools given in Section 4 are used; these include bias, variances for parametric components, MSEs for nonparametric components, variance of the model, and REs for the all models. Results are given following tables and figures

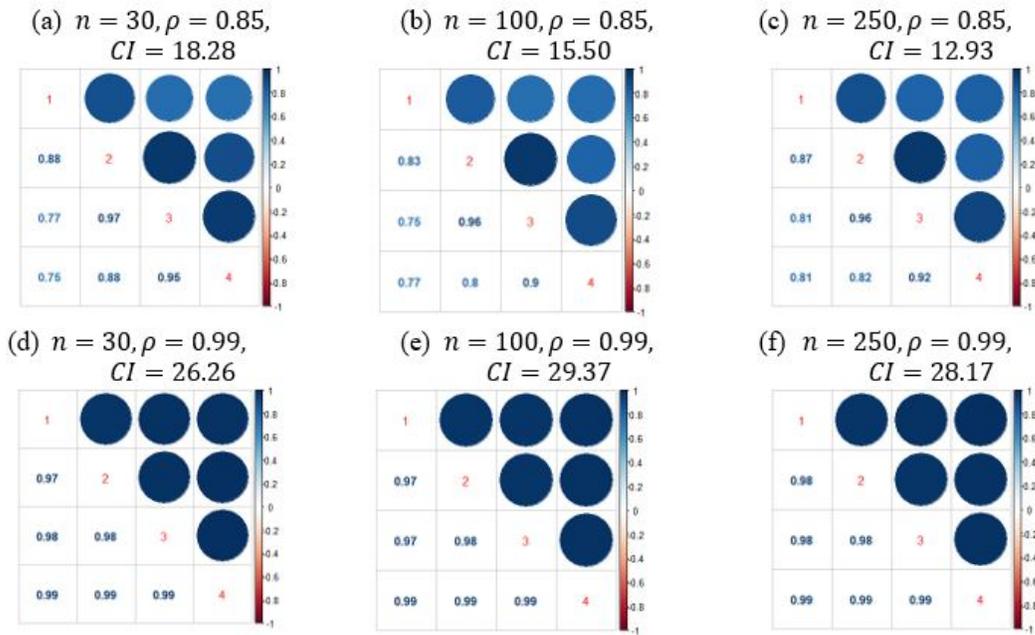


Figure 1: Correlation plots for correlation levels  $\rho = 0.85$  and  $\rho = 0.99$  for three sample sizes.

Figure 1 depicts correlation values between covariates that are possessed by the parametric component of the model. Also, condition index ( $CI$ ) values are added to each panel of the figure. As can be seen, all  $CI$  values are larger than 10, which indicates serious multicollinearity. It is ensured that the two correlation levels ( $\rho = 0.85, \rho = 0.99$ ) for explanatory variables are seen clearly in the sixth panel of the figure.

Table 1: Outcomes for the parametric components when  $\rho = 0.85$ 

		$\rho = 0.85$							
		0.3				1			
$\sigma_\varepsilon$	Method	SDT		kNN		SDT		kNN	
CL	n	Var( $\beta$ )	Bias( $\beta$ )						
5%	30	<b>1.5299</b>	<b>0.2340</b>	2.1735	0.2346	<b>1.9954</b>	<b>0.2165</b>	2.5834	0.2281
	100	<b>0.8309</b>	0.2013	1.2315	<b>0.1732</b>	1.2610	0.2009	<b>0.8845</b>	<b>0.1763</b>
	250	0.4075	0.1555	<b>0.3626</b>	<b>0.1277</b>	0.3794	0.1657	<b>0.3156</b>	<b>0.1325</b>
35%	30	1.2576	0.2813	<b>1.2207</b>	<b>0.2641</b>	2.1904	0.2825	<b>1.7777</b>	<b>0.2636</b>
	100	<b>0.2500</b>	0.2580	0.4861	<b>0.2027</b>	<b>0.2765</b>	0.2582	0.7216	<b>0.2046</b>
	250	0.2900	0.1989	<b>0.2525</b>	<b>0.1445</b>	0.2885	0.1841	<b>0.2435</b>	<b>0.1421</b>

Performance scores of methods on estimating parametric component of the model are presented Tables 1 and 2 for  $\rho = 0.85$  and  $\rho = 0.99$ , respectively. The best scores are indicated with bold font. It can be observed that the quality of estimations gets worse when censoring level and error variance increase, which is an expected situation. This case is also true of the correlation level. It is obvious that correlation level affects the performance of estimates, which can be monitored by Figures 3 and 4.

An interesting result from these tables is that it can be concluded that the kNN based model gives less bias than the synthetic data transformation (SDT) based model. Although both methods provide close scores in terms of  $Var(\beta)$ , kNN dominates the SDT method in  $Bias(\beta)$ , especially for high censoring levels. The reason for that can be explained by the nature of SDT, which causes spatial variation in the data. Because kNN works with data points individually, it manipulates data structure less than the SDT method.

Figures 2 and 3 are drawn to emphasize the bias of the regression coefficients. S1 and K1 in the x-axes in each panel represent SDT and kNN based biases for  $n = 30$ . Similarly, S2 and K2 denote biases for  $n = 100$ , S3 and K3 for  $n = 250$ . It can be seen from an examination of these figures that the boxplots are similar for both methods. In some cases, such as  $\sigma = 0.3$ , the SDT based estimated values of  $\beta$  are slightly better than those of kNN, but when variation in the generated data increases, creates better estimates. This can be seen in bottom panels of both figures. It can therefore be said that Figures 3 and 4 verify the results given in Tables 1 and 2.

Table 2: Outcomes for the parametric components when  $\rho = 0.99$

		$\rho = 0.99$							
$\sigma_\varepsilon$		0.3				1			
Method		SDT		kNN		SDT		kNN	
CL	n	Var( $\beta$ )	Bias( $\beta$ )						
	30	<b>9.0435</b>	<b>0.2978</b>	10.6518	0.315	<b>10.2537</b>	<b>0.2974</b>	11.1486	0.3018
5%	100	<b>4.1748</b>	0.2255	4.9727	<b>0.1889</b>	5.8733	0.2143	<b>5.0571</b>	<b>0.1799</b>
	250	2.1894	0.1538	<b>1.9546</b>	<b>0.1169</b>	<b>2.6385</b>	0.1524	2.7314	<b>0.1137</b>
	30	<b>15.0880</b>	<b>0.3115</b>	15.6767	0.3461	<b>11.7627</b>	0.3155	12.5868	<b>0.3144</b>
35%	100	6.8527	0.3042	<b>6.2235</b>	<b>0.2010</b>	6.7189	0.3074	<b>6.7090</b>	<b>0.2060</b>
	250	4.1108	0.2050	<b>4.0858</b>	<b>0.1274</b>	5.5155	0.2044	<b>3.1370</b>	<b>0.1279</b>

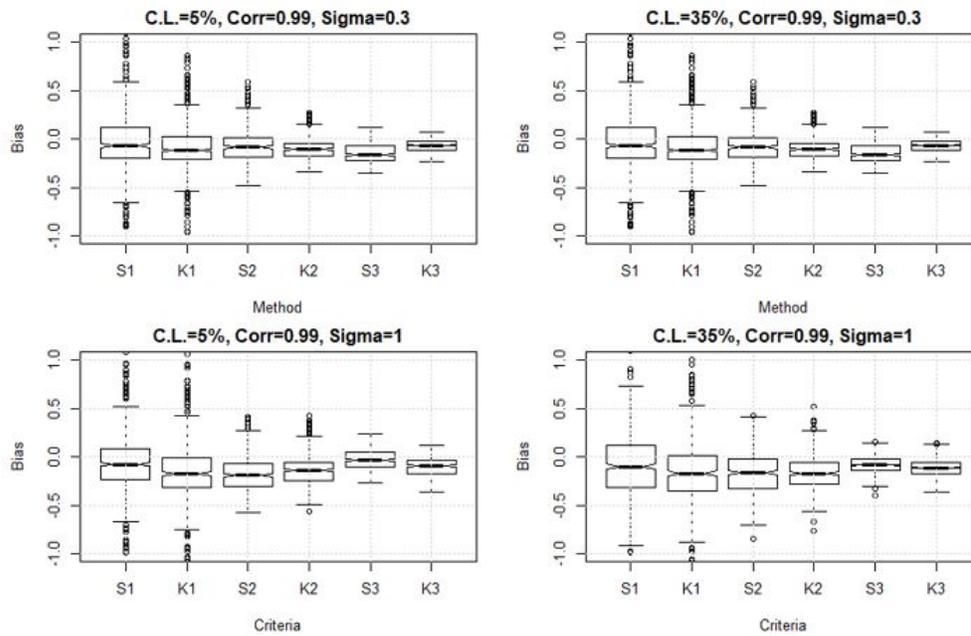


Figure 2: Boxplots for bias of regression coefficients

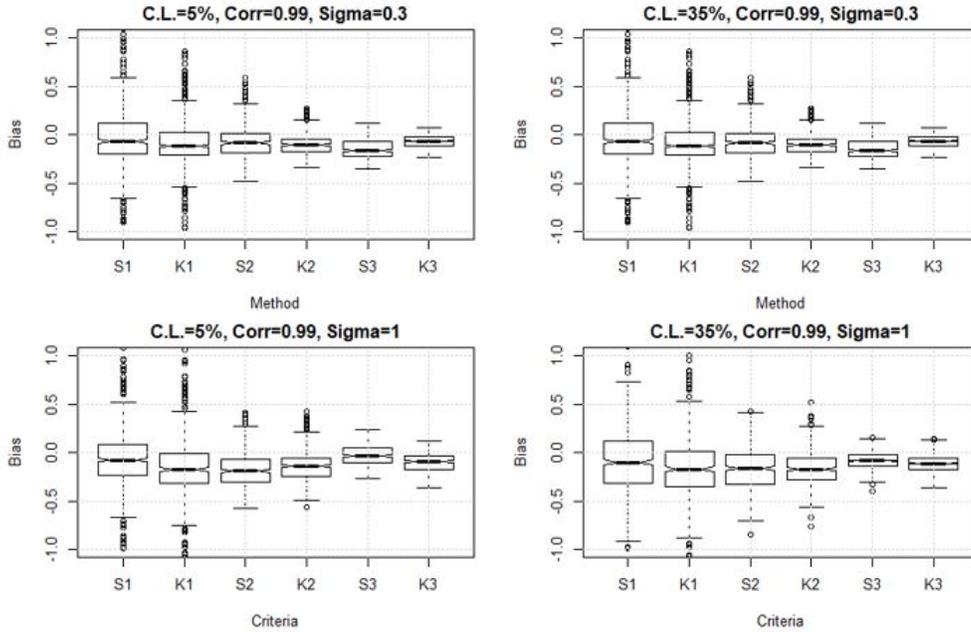


Figure 3: Boxplots for bias of regression coefficients

In Table 3, outcomes for the nonparametric component  $g(z)$  are given. As noted above, MSE is commonly used as an evaluation metric for  $g(z)$ . From obtained MSEs in the Table 3, it can be seen that in most combinations, the scores are the same or close to each other. This can be confirmed from Figures 4-5. However, for different correlation levels and error variances, the outcomes change in favor of kNN based estimations. Differences between the methods become clear under hard conditions, such as heavy censoring or high correlation levels. From this, it can be said that although both methods provide satisfying performances on estimation of the nonparametric component, the kNN based model is more resistant to difficult conditions.

Figures 4 and 5 depict the fitted curves of SDT and kNN based estimators. Note that, because there are too many combinations to show, only some selections are presented. Figure 4 shows the estimated nonparametric functions for low correlation level  $\rho = 0.85$  and Figure 5 is obtained from  $\rho = 0.99$ .

Table 3: Outcomes for the nonparametric components

		$\rho = 0.85$				$\rho = 0.99$			
$\sigma$		0.3		1		0.3		1	
CL	n	SDT	kNN	SDT	kNN	SDT	kNN	SDT	kNN
5%	30	0.0735	0.0735	0.0736	0.0736	0.0818	0.0818	0.0866	<b>0.0865</b>
	100	0.0718	0.0718	0.0768	0.0768	0.0853	0.0853	0.0903	0.0903
	250	0.0777	<b>0.0776</b>	<b>0.0771</b>	0.0776	0.0864	<b>0.0862</b>	0.0914	0.0914
35%	30	0.0920	<b>0.0915</b>	0.1052	<b>0.1046</b>	0.0981	<b>0.0973</b>	0.1131	<b>0.1123</b>
	100	0.0960	<b>0.0957</b>	0.1097	<b>0.1094</b>	0.1024	<b>0.1020</b>	0.1182	<b>0.1178</b>
	250	0.0971	<b>0.0970</b>	0.1110	<b>0.1108</b>	0.1036	<b>0.1034</b>	0.1195	<b>0.1193</b>

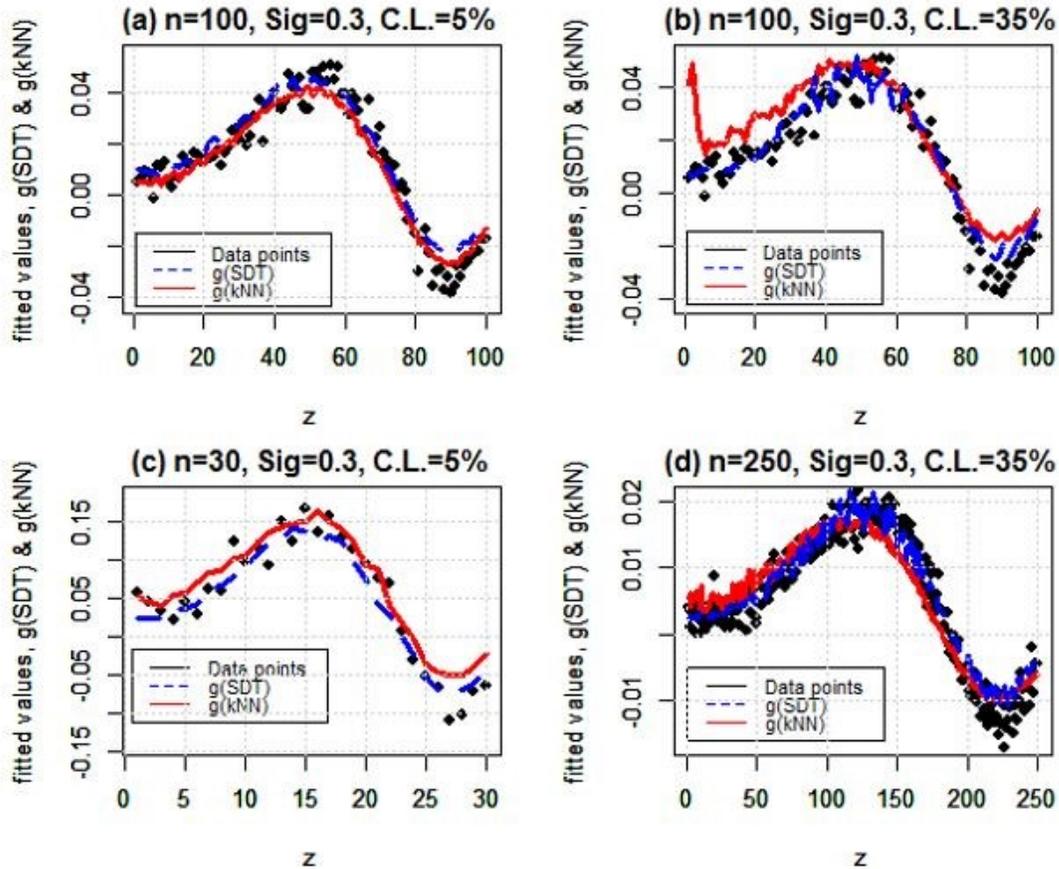


Figure 4: Fitted curves for synthetic data transformation (blue line), kNN imputation technique (red line) and real function (black line) when  $\rho = 0.85$ .

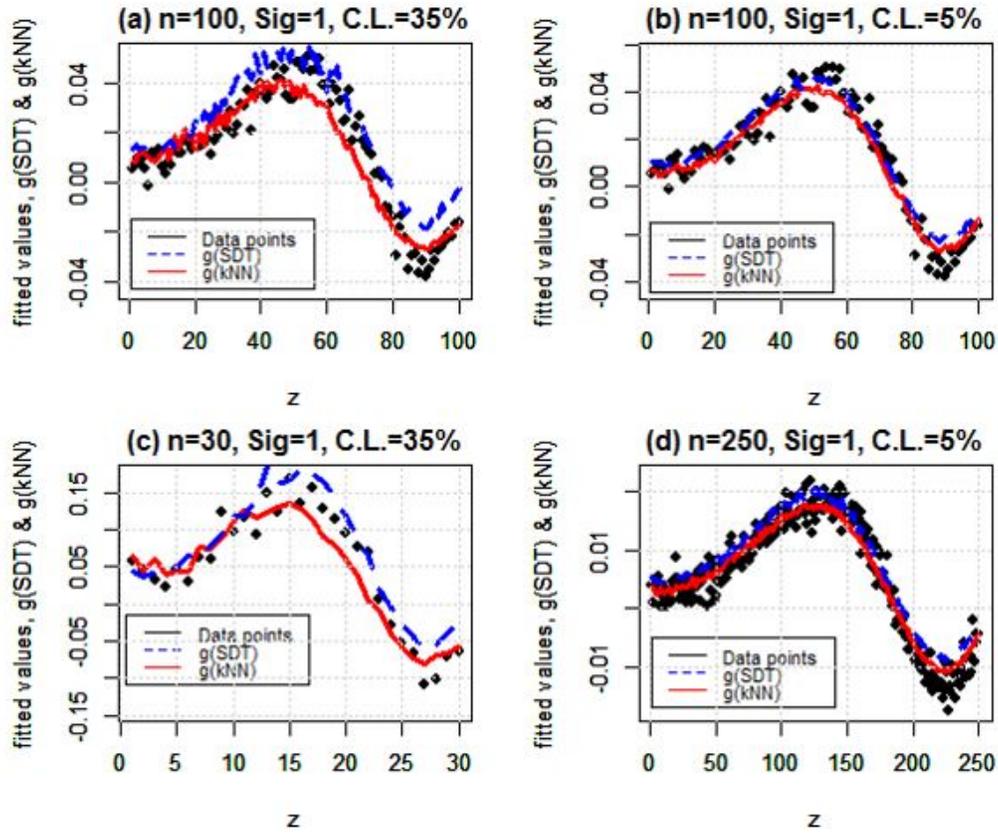


Figure 5: Fitted curves for synthetic data transformation (blue line), kNN imputation technique (red line) and real function (black line) when  $\rho = 0.99$ .

In the top-right panel of Figure 4, one can observe that the "weird" curve obtained by kNN seems to contradict the simulation experiments, because kNN has a better performance compared to SDT. This is because at this point, the nonparametric nature of the kNN method comes into play. Because of kNN's nature, it has the potential to give some outlier results. The plot given in top-right panel of Figure 4 can be counted as an example of such an outlier. The remaining plots in Figures 4-5 support the results in Table 3, depicting close-fitted curves and better kNN results under difficult conditions, as can be seen in the bottom-left of Figure 5.

Table 4: The REs and estimated model variances

$\rho = 0.85$									
$\sigma$		0.3				1			
Method		SDT		kNN		SDT		kNN	
CL	n	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE
	30	0.9124	1.4067	<b>0.6486</b>	<b>0.7109</b>	<b>0.6159</b>	<b>0.9583</b>	0.6427	1.0435
5%	100	0.9020	1.4277	<b>0.6318</b>	<b>0.7004</b>	0.6420	1.0436	<b>0.6152</b>	<b>0.9583</b>
	250	0.9059	1.4282	<b>0.6343</b>	<b>0.7002</b>	0.6503	1.0448	<b>0.6224</b>	<b>0.9571</b>
	30	0.8906	1.4139	<b>0.6299</b>	<b>0.7073</b>	0.8840	1.4023	<b>0.6304</b>	<b>0.7131</b>
35%	100	0.8895	1.4347	<b>0.6200</b>	<b>0.6970</b>	0.8937	1.4391	<b>0.6210</b>	<b>0.6949</b>
	250	0.9037	1.4471	<b>0.6245</b>	<b>0.6910</b>	0.9001	1.4452	<b>0.6228</b>	<b>0.6919</b>
$\rho = 0.99$									
$\sigma$		0.3				1			
Method		SDT		kNN		SDT		kNN	
CL	n	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE	$\sigma_\varepsilon^2$	RE
	30	0.6467	1.0534	<b>0.6139</b>	<b>0.9493</b>	<b>0.6196</b>	<b>0.9487</b>	0.6531	1.0541
5%	100	0.6459	1.0595	<b>0.6096</b>	<b>0.9438</b>	<b>0.6065</b>	<b>0.9888</b>	0.6134	1.0114
	250	<b>0.5817</b>	<b>0.9632</b>	0.6039	1.0382	0.5743	1.0058	<b>0.5710</b>	<b>0.9943</b>
	30	0.8971	1.3913	<b>0.6448</b>	<b>0.7188</b>	0.9124	1.2188	<b>0.7486</b>	<b>0.8205</b>
35%	100	0.8763	1.3927	<b>0.6292</b>	<b>0.7180</b>	0.9020	1.2326	<b>0.7318</b>	<b>0.8113</b>
	250	0.8335	1.3763	<b>0.6056</b>	<b>0.7266</b>	0.8959	1.2369	<b>0.7243</b>	<b>0.8085</b>

Finally, Table 4 involves the model outcomes that are REs and estimated model variances. Note that because the nonparametric results are close to each other, and kNN has better scores in terms of parametric component of the model, in the model performance, the kNN method is indirectly given better model scores. Here, it can once more be seen that kNN reveals its difference for high correlation levels and high censoring levels. Otherwise, SDT has an acceptable performance for the estimation of a semiparametric model under censored multicollinear data. For this study, the kNN method demonstrates important success on modeling multicollinear censored data, which can be clearly seen in Table 4.

From the information obtained from the simulation experiments, it can be said that kernel ridge estimators based on SDT and kNN methods can handle multicollinear datasets. In addition, this study shows that kNN imputation gives more satisfying

results than SDT, which is an important conclusion. Details about the results are given in Section 6.

## 6 Conclusion

This paper considers the estimation of a partially linear model under multicollinear censored data. The aim of the paper is realized using a kernel-type ridge estimator based on two different censorship solution techniques: synthetic data transformation (SDT) and kNN imputation, respectively. This study therefore adds two main contributions to the literature. To show the behaviors of the introduced estimators, a detailed simulation study was designed and performed. From the results obtained from this simulation study (presented in Section 5), concluding remarks are given as follows:

- Tables 1 and 2 represent the outcomes of the parametric component of the model. It can be said that both methods, SDT and kNN, have satisfying results for estimating semiparametric models under data irregularities. On the other hand, the kNN based kernel ridge-type estimator performs better than the SDT based model, especially under heavy censoring levels and high correlation levels. This is because SDT manipulates the data set.
- Figures 2 and 3 show boxplots for the biases of the regression coefficients and confirm the results of Tables 1 and 2. The figures also demonstrate the effects of sample size, censoring level, correlation level, and error variance on the estimating of regression coefficients.
- Table 3 includes MSE values for the nonparametric component of the model for both estimators, kNN and SDT. The simulation study shows that, on the estimation of nonparametric function, the methods give highly similar results except for under high censoring and high correlation levels. The kNN method is more resistant to irregularities than SDT. Figures 4 and 5 show the fitted curves and clearly illustrate the inferences which can be drawn from Table 3.
- Finally, Table 4 shows the performance scores for model variance and relative efficiencies (RE) for the estimated models. The kNN based model has better estimates than SDT, although it should be noted that, for low censoring and correlation levels, SDT has also satisfying results. Thus, both estimators can be used under appropriate conditions.

## Acknowledgement

The authors are thankful to the guest editor and anonymous reviewer for constructive comments. Prof. Dr. S. Ejaz Ahmed's research is supported by the Natural Sciences and the Engineering Research Council of Canada (NSERC).

## References

- Ahmed, S. E., Aydın, D., and Yilmaz, E. (2020), Nonparametric regression estimates based on imputation techniques for right-censored data. *Proceedings of the Thirteenth International Conference on Management Science and Engineering Management, ICMSEM 2019, Advances in Intelligent Systems and Computing*, Vol **1001**, Springer, Cham.
- Ahmed, S. E. (2014), *Penalty, Shrinkage and Pretest Strategies: Variable Selection and Estimation*. Springer, New York, <https://www.springer.com/gp/book/9783319031484>.
- Batista, G., and Monard, M. (2002), An analysis of four missing data treatment methods for supervised learning, *Applied Artificial Intelligence*, **17**, 519-533.
- Hoerl, A. E., and Kennard, R. W. (1970), Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, **12**(1), 55-67.
- Kaplan, E. L., and Meier, P. (1958), Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, **53**(282), 457-481.
- Koul, H., Susarla, V., and Van Ryzin, J. (1981), Regression analysis with randomly right censored data. *The Annals of Statistics*, **9**, 1276-88.
- Liang, H. (2006), Estimation in partially linear models and numerical comparisons. *Computational Statistics and Data Analysis*, **50**(3), 675-685.
- Liang, H., and Zhou, Y. (2008), Semiparametric Inference for ROC curves with censoring. *Scandinavian Journal of Statistics*, **35**(2), 212-227.
- Miller, R. G. (1976), Least squares regression with censored data. *Biometrika*, **63**, 449-64.
- Orbe, J., Ferreira, E., and Nunez-Anton, V. (2003), Censored Partial Regression. *Biostatistics*, **4**(1), 109-121.
- Ruppert, D., Wand, M. P., and Carroll, R. J. (2003), Semiparametric regression. *Cambridge University Press*, New York.

- Schimek, M. G. (2000), Smoothing and Regression: Approaches, Computation and Application. *John Wiley & Sons*, Hoboken, NJ.
- Shim, J. Y. (2005), Censored kernel ridge regression. *Journal of the Korean Data and Information Science Society*, **16**(4), 1045-1052.
- Stute, W. (1993), Consistent Estimation Under Random Censorship When Covariables Are Present. *Journal Of Multivariate Analysis*, **4**, 89-103.
- Yenduri, S., and Iyengar S. S. (2007), Performance evaluation of imputation methods for incomplete datasets. *International Journal of Software Engineering and Knowledge Engineering*, **17**, 127-152.
- Yüzbaşı, B., Arashi, M., and Ahmed, S. E. (2017), Shrinkage estimation strategies in generalized ridge regression models under low/high-dimension regime. *International Statistical Review*, **88**(1), 229-251.

## Appendix

### A1. Proof of Theorem 3.1

As mentioned before,  $\mathbf{W}_h$  is the kernel smoother matrix and  $\tilde{\mathbf{X}}$  in (3.14) are the variables obtained by using  $\mathbf{W}_h$ , as mentioned in Speckman (1988). From that, let us consider the weighted minimization problem to obtain (3.14)-(3.16) as follows:

$$\begin{aligned} L_{\hat{G}}(\boldsymbol{\beta}) &= \sum_{i=1}^n (\varepsilon_{i\hat{G}})^2 = \sum_i (\tilde{t}_{i\hat{G}} - \tilde{\mathbf{X}}_i \boldsymbol{\beta})^2 + \sum_{j=1}^n (0 - k\beta_j)^2 \\ &= (\tilde{\mathbf{t}}_{\hat{G}} - \tilde{\mathbf{X}} \boldsymbol{\beta})^T (\tilde{\mathbf{t}}_{\hat{G}} - \tilde{\mathbf{X}} \boldsymbol{\beta}) + k \|\mathbf{0} - \boldsymbol{\beta}\|^2. \end{aligned} \quad (\text{A1.1})$$

In order to minimize (A1.1), augmented data sets are used. This can be shown as

$$\mathbf{X}_A = \begin{pmatrix} \tilde{\mathbf{X}}_{n \times p} \\ \sqrt{k} \mathbf{I}_p \end{pmatrix} \text{ and } \mathbf{t}_A = \begin{pmatrix} \tilde{\mathbf{t}}_{\hat{G}}_{n \times 1} \\ \mathbf{0}_p \end{pmatrix}, \quad (\text{A1.2})$$

where  $\mathbf{I}_p$  is the  $(p \times p)$  dimensional identity matrix and  $\mathbf{0}_p$  is the  $(n \times 1)$  zero matrix. For details of (A1.2), see Aydın et al. (2018). It should be differentiated with respect to  $\boldsymbol{\beta}$  and set the equation to zero, which is given by

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} L_{\hat{G}}(\boldsymbol{\beta}) &= -2(\tilde{\mathbf{X}} \tilde{\mathbf{t}}_{\hat{G}}) + 2\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\beta} = \mathbf{0}, \\ \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\beta} &= \tilde{\mathbf{X}}_{\hat{G}}, \\ \hat{\boldsymbol{\beta}}_{\hat{G}} &= (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k \mathbf{I}_p)^{-1} \tilde{\mathbf{X}}^T \mathbf{t}_{\hat{G}}. \end{aligned} \quad (\text{A1.3})$$

From (A1.3),  $\hat{\boldsymbol{\beta}}_{\hat{G}}$  is obtained as follows easily as in equation (3.16):

$$\hat{\boldsymbol{\beta}}_{\hat{G}} = \mathbf{W}_h (\tilde{\mathbf{t}}_{\hat{G}} - \mathbf{X}^T \hat{\boldsymbol{\beta}}_{\hat{G}}). \quad (\text{A1.4})$$

### Proof of Theorem 3.2

Similar to Appendix A1, replacing synthetic response variable  $\tilde{t}_{\hat{G}}$  by the kNN based response variable  $\mathbf{y}^k$  minimization criterion is written as follows:

$$\begin{aligned} L_{kNN}(\boldsymbol{\beta}) &= \sum_{i=1}^n (\varepsilon_i)^2 = \sum_i (\tilde{y}_i^k - \tilde{\mathbf{X}}_i \boldsymbol{\beta})^2 + \sum_{j=1}^n (0 - k\beta_j)^2 \\ &= (\tilde{\mathbf{y}}^k - \tilde{\mathbf{X}} \boldsymbol{\beta})^T (\tilde{\mathbf{y}}^k - \tilde{\mathbf{X}} \boldsymbol{\beta}) + k \|\mathbf{0} - \boldsymbol{\beta}\|^2. \end{aligned} \quad (\text{A2.1})$$

Then, as in A1, augmented data sets are given by

$$\mathbf{X}_A = \begin{pmatrix} \tilde{\mathbf{X}}_{n \times p} \\ \sqrt{k} \mathbf{I}_p \end{pmatrix} \text{ and } \mathbf{t}_A = \begin{pmatrix} \tilde{\mathbf{y}}_{n \times 1}^k \\ \mathbf{0}_p \end{pmatrix}, \quad (\text{A2.2})$$

where  $\mathbf{I}_p$  is the  $(p \times p)$  dimensional identity matrix and  $\mathbf{0}_p$  is the  $(n \times 1)$  zero matrix. According to (A2.2), derivation of  $\boldsymbol{\beta}$  is given as

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} L_{kNN}(\boldsymbol{\beta}) &= -2(\tilde{\mathbf{X}} \tilde{\mathbf{y}}^k) + 2\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\beta} = \mathbf{0}, \\ \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \boldsymbol{\beta} &= \tilde{\mathbf{X}} \tilde{\mathbf{y}}^k, \\ \hat{\boldsymbol{\beta}}_k &= \left( \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + k \mathbf{I}_p \right)^{-1} \tilde{\mathbf{X}}^T \tilde{\mathbf{y}}^k. \end{aligned} \quad (\text{A2.3})$$

Thus equation (3.17) is derived. From (A2.3),  $\hat{\boldsymbol{g}}_{\hat{\mathbf{C}}}$  is obtained as follows easily as in equation (3.18):

$$\hat{\boldsymbol{g}}_{\hat{\mathbf{C}}} = \mathbf{W}_h \left( \tilde{\mathbf{y}}^k - \tilde{\mathbf{X}}^T \hat{\boldsymbol{\beta}}_k \right). \quad (\text{A2.4})$$