

Estimation in monotone single-index models

Piet Groeneboom¹ | Kim Hendrickx²

¹Delft Institute of Applied Mathematics, Delft University of Technology, Delft, The Netherlands

²I-BioStat, Hasselt University, Hasselt, Belgium

Correspondence

Piet Groeneboom, Delft Institute of Applied Mathematics, Delft University of Technology, Delft, The Netherlands.
Email: P.Groeneboom@tudelft.nl

Funding information

Research Foundation Flanders (FWO), Grant/Award Number: 11W7315N; IAP Research Network P7/06 of the Belgian State (Belgian Science Policy); VSC-Flemish Supercomputer Center, funded by the Hercules Foundation and the Flemish Government-department EWI

Single-index models are popular regression models that are more flexible than linear models and still maintain more structure than purely nonparametric models. We consider the problem of estimating the regression parameters under a monotonicity constraint on the unknown link function. In contrast to the standard approach of using smoothing techniques, we review different “non-smooth” estimators that avoid the difficult smoothing parameter selection. For about 30 years, one has had the conjecture that the profile least squares estimator is an \sqrt{n} -consistent estimator of the regression parameter, but the only non-smooth argmin/argmax estimators that are actually known to achieve this \sqrt{n} -rate are not based on the nonparametric least squares estimator of the link function. However, solving a score equation corresponding to the least squares approach results in \sqrt{n} -consistent estimators. We illustrate the good behavior of the score approach via simulations. The connection with the binary choice and current status linear regression models is also discussed.

KEYWORDS

least squares, monotone link function, single-index model

1 | THE SINGLE-INDEX MODEL

Suppose that Y is a response variable and that $\mathbf{X} = (X_1, \dots, X_d)^T$ is a d -dimensional covariate ($d \geq 1$). The semiparametric single-index model is given by

$$Y = \psi_0(\boldsymbol{\alpha}_0^T \mathbf{X}) + \varepsilon, \quad (1)$$

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2018 The Authors. *Statistica Neerlandica* published by John Wiley & Sons Ltd on behalf of VVS.

where $\varepsilon \sim F_0$ is a random error term with $\mathbb{E}(\varepsilon|\mathbf{X}) = 0$ almost surely, ψ_0 is an unknown function, and α_0 is an unknown regression parameter in \mathbb{R}^d . These models are more flexible than standard linear regression models and have, on the other hand, more structure than completely nonparametric models.

In this paper, we impose a monotonicity constraint on the link function ψ_0 . Shape-constrained inferences have gained a lot of popularity in recent years. Several statistical applications are based on imposing constraints that occur from the problem under study and need algorithms for shape-constrained regression or density estimation. Monotonicity, convexity, or concavity constraints arise naturally with consumption or production functions, growth curves, and dose response models. A summary of nonparametric estimation under shape constraints can be found in Groeneboom and Jongbloed (2014).

Model (1), when ψ_0 is an unknown monotone function, is also known as the monotone single-index model, and many econometric models as well as censored regression models and various duration models fit into this framework. We next describe how the binary choice model and the current status linear regression model are special cases of the monotone single-index model satisfying $\mathbb{E}(Y|\mathbf{X}) = \psi_0(\alpha_0^T \mathbf{X})$.

The binary choice model. A widely used econometric model is the binary choice model, which is used to describe a choice probability based on one or more covariates. The model is given by

$$Y = \begin{cases} 1, & \text{if } \alpha_0^T \mathbf{X} \geq \varepsilon \\ 0, & \text{else,} \end{cases} \quad (2)$$

where $\alpha_0^T \mathbf{X}$ presents the utility score and ε , assumed to be independent of \mathbf{X} , is the disturbance term. The model can be used to predict the probability that a person decides to consume a certain good based on the characteristics of the person. The model is a special case of the single-index model (1) with ψ_0 equal to the (unknown) distribution function F_0 of ε , since

$$\mathbb{E}\{Y|\mathbf{X}\} = P(Y = 1|\mathbf{X}) = P(\varepsilon \leq \alpha_0^T \mathbf{X}) = F_0(\alpha_0^T \mathbf{X}).$$

The current status linear regression model. A frequently encountered problem in regression analysis is that a variable of interest (which can be thought of as an event time) is not observed directly but only known to lie before or after some random censoring time. This type of censored data is known as current status data and arises commonly in reliability and survival studies especially when testing is destructive. More formally, instead of observing the vector (Y, \mathbf{X}) as in (1), a vector (T, Δ, \mathbf{X}) is observed, where T is a censoring variable and $\Delta = 1_{Y \leq T}$. Since

$$\mathbb{E}\{\Delta|T, \mathbf{X}\} = P(\Delta = 1|T, \mathbf{X}) = F_0(T - \alpha_0^T \mathbf{X}),$$

the current status linear regression model is a special case of the single-index model (1) with response $\tilde{Y} = \Delta$, covariate vector $\tilde{\mathbf{X}} = (T, \mathbf{X})^T$, and $\tilde{\alpha}_0 = (1, \alpha_0^T)$.

The models are prototypes of semiparametric models, where one has a nonparametric component given by the “link function” ψ_0 , which is a distribution function (or 1 minus a distribution function) in the current status regression model, and a parametric part. The difficulty of the model is that the parametric part is “inside” the nonparametric part; one has to bypass the nonparametric function ψ_0 , which cannot be estimated at rate \sqrt{n} , to get to the parametric part.

This is very different for the well-known Cox proportional hazards model for current status data. In this case, the log likelihood is of the form

$$\sum_{i=1}^n \left\{ \Delta_i \log \left(1 - \exp \left\{ -\Lambda(T_i) e^{\alpha_0^T X_i} \right\} \right) - (1 - \Delta_i) \Lambda(T_i) e^{\alpha_0^T X_i} \right\},$$

where Λ is the baseline cumulative hazard function. Now, α_0 does not appear in the argument of a function F , which is not \sqrt{n} estimable, and we can estimate Λ and α_0 separately. In this case, it was shown in Huang (1996) that one can use the nonparametric maximum likelihood estimator of Λ and then use profile likelihood to estimate α_0 efficiently at rate \sqrt{n} . However, for the ordinary current status regression model, it is still unknown whether a similar estimation method gives an \sqrt{n} -consistent estimate of α_0 .

1.1 | Identifiability

Identification of the single-index regression parameter α_0 (up to a scalar constant) has been discussed in Ichimura (1993) in terms of the distribution of the regressors \mathbf{X} . Without any further restrictions, the parameter vector (α_0, ψ_0) , however, cannot be estimated in the single-index model. This can be seen as follows. Take $a, b \in \mathbb{R}$ and let ψ^* be the function defined by the relationship $\psi^*(a + bt) = \psi_0(t)$ for all t in the support of $\alpha_0^T \mathbf{X}$, then

$$\mathbb{E}(Y|\mathbf{X}) = \psi^*(a + b\alpha_0^T \mathbf{X}).$$

Even if the distribution of (Y, \mathbf{X}) is known, the above model cannot be distinguished from model (1) unless restrictions on location a and scale b are imposed. Location normalization can be imposed by requiring that all components of \mathbf{X} have a nondegenerate distribution. A reparametrization of the parameter space to the set

$$\{\alpha \in \mathbb{R}^d : \alpha_1 = 1\} \quad \text{or} \quad \{\alpha \in \mathbb{R}^d : \|\alpha_0\| = 1, \alpha_{01} \geq 0\},$$

where $\|\cdot\|$ denotes the Euclidean norm and α_{01} is the first component of α_0 , ensures scale identification of the model. The first parametrization is used in Sherman (1993), whereas examples of the second parametrization are found in Härdle, Hall, and Ichimura (1993) and Hristache, Juditsky, and Spokoiny (2001), among others. Note that in the special case of the current status linear regression model, the first component of the covariate vector corresponds to the censoring variable T with a coefficient equal to 1, and therefore, the current status regression model is identified without further restrictions on the parameter space.

1.2 | Efficient information

The log likelihood for the model of one observation is given by

$$\ell_{\alpha, \psi}(\mathbf{x}, y) \stackrel{\text{def}}{=} \log \{f_{\varepsilon|\mathbf{X}}(y - \psi(\alpha^T \mathbf{x})) f_{\mathbf{X}}(\mathbf{x})\},$$

where $f_{\varepsilon|\mathbf{X}}$ is the conditional density of ε given $\mathbf{X} = \mathbf{x}$ and $f_{\mathbf{X}}$ is the density of \mathbf{X} .

The partial derivative w.r.t. α of $\ell_{\alpha, \psi}$ is given by

$$\frac{\partial}{\partial \alpha} \ell_{\alpha, \psi}(\mathbf{x}, y) = \frac{\mathbf{x} \psi'(\alpha^T \mathbf{x}) f'_{\varepsilon|\mathbf{X}}(y - \psi(\alpha^T \mathbf{x}))}{f_{\varepsilon|\mathbf{X}}(y - \psi(\alpha^T \mathbf{x}))}.$$

Let $\{\psi_\eta : \eta \in (-1, 1)\}$ be a path in the collection $\{\psi : \psi \text{ is increasing}\}$, differentiable w.r.t. η at $\eta = 0$, and suppose

$$\psi_\eta = \psi \quad \text{for } \eta = 0$$

and

$$\left. \frac{\partial}{\partial \eta} \psi_\eta(t) \right|_{\eta=0} = a(t).$$

Then

$$\frac{\partial}{\partial \eta} \ell_{\alpha, \psi_\eta}(\mathbf{x}, y) \Big|_{\eta=0} = \frac{a(\alpha^T \mathbf{x}) f'_{\varepsilon|\mathbf{X}}(y - \psi(\alpha^T \mathbf{x}))}{f_{\varepsilon|\mathbf{X}}(y - \psi(\alpha^T \mathbf{x}))}.$$

To obtain the efficient score function, we must solve

$$\mathbb{E} \left[\left\{ \frac{\mathbf{X} \psi'(\alpha^T \mathbf{X}) f'_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))}{f_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))} - \frac{a_*(\alpha^T \mathbf{X}) f'_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))}{f_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))} \right\} \frac{a(\alpha^T \mathbf{X}) f'_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))}{f_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))} \right] = 0, \quad (3)$$

for an \mathbb{R}^d -valued function a_* , where $a_*, a \in L^0_2(F)^d$ (see, e.g., Huang, 1996, p. 558, for similar computations with an \mathbb{R} -valued function a_*). This amounts to solving in a_*

$$\mathbb{E} \left[\frac{\{\mathbf{X} \psi'(\alpha^T \mathbf{X}) - a_*(\alpha^T \mathbf{X})\} f'_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))^2}{f_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))^2} a(\alpha^T \mathbf{X}) \right] = 0.$$

The efficient variance for α in the single-index model is derived in Newey and Stoker (1993), Delecroix, Härdle, and Hristache (2003), and Kuchibhotla and Patra (2017), among others. For the general case, we get that the efficient score function is given by

$$\tilde{\ell}_{\alpha, \psi}(\mathbf{x}, y) = \frac{y - \psi(\alpha^T \mathbf{x})}{\sigma^2(\mathbf{x})} \psi'(\alpha^T \mathbf{x}) \left\{ \mathbf{x} - \frac{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) \mathbf{X} | \alpha^T \mathbf{X} = \alpha^T \mathbf{x} \}}{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) | \alpha^T \mathbf{X} = \alpha^T \mathbf{x} \}} \right\}, \quad (4)$$

where $\sigma^2(\cdot) = \mathbb{E}(\varepsilon^2 | \mathbf{X} = \cdot)$. We illustrate the derivation of this efficient score function in case that $\varepsilon | \mathbf{X} \sim N(0, \sigma^2(\mathbf{X}))$. We can write

$$\begin{aligned} & \mathbb{E} \left[\frac{\{\mathbf{X} \psi'(\alpha^T \mathbf{X}) - a_*(\alpha^T \mathbf{X})\} f'_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))^2}{f_{\varepsilon|\mathbf{X}}(Y - \psi(\alpha^T \mathbf{X}))^2} a(\alpha^T \mathbf{X}) \right] \\ &= \mathbb{E} \left[\{\mathbf{X} \psi'(\alpha^T \mathbf{X}) - a_*(\alpha^T \mathbf{X})\} \frac{\{y - \psi(\alpha^T \mathbf{X})\}^2}{\sigma^4(\mathbf{X})} a(\alpha^T \mathbf{X}) \right] \\ &= \mathbb{E} \left[\mathbb{E} \left\{ \frac{\{\mathbf{X} \psi'(\alpha^T \mathbf{X}) - a_*(\alpha^T \mathbf{X})\}}{\sigma^2(\mathbf{X})} \Big| \alpha^T \mathbf{X} \right\} a(\alpha^T \mathbf{X}) \right]. \end{aligned}$$

Note that

$$\begin{aligned} & \mathbb{E} \left\{ \frac{\{\mathbf{X} \psi'(\alpha^T \mathbf{X}) - a_*(\alpha^T \mathbf{X})\}}{\sigma^2(\mathbf{X})} \Big| \alpha^T \mathbf{X} \right\} \\ &= \psi'(\alpha^T \mathbf{X}) \mathbb{E} \left\{ \sigma^{-2}(\mathbf{X}) \mathbf{X} \Big| \alpha^T \mathbf{X} \right\} - a_*(\alpha^T \mathbf{X}) \mathbb{E} \left\{ \sigma^{-2}(\mathbf{X}) \Big| \alpha^T \mathbf{X} \right\}. \end{aligned}$$

This means that (3) is solved for

$$a_*(u) = \psi'(u) \frac{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) \mathbf{X} | \alpha^T \mathbf{X} = u \}}{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) | \alpha^T \mathbf{X} = u \}}.$$

We conclude that the efficient score function for the semiparametric single-index model if $\varepsilon | \mathbf{X} \sim N(0, \sigma^2(\mathbf{X}))$ is indeed given by

$$\tilde{\ell}_{\alpha, \psi}(\mathbf{x}, y) = \frac{y - \psi(\alpha^T \mathbf{x})}{\sigma^2(\mathbf{x})} \psi'(\alpha^T \mathbf{x}) \left\{ \mathbf{x} - \frac{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) \mathbf{X} | \alpha^T \mathbf{X} = \alpha^T \mathbf{x} \}}{\mathbb{E} \{ \sigma^{-2}(\mathbf{X}) | \alpha^T \mathbf{X} = \alpha^T \mathbf{x} \}} \right\}.$$

Efficiency calculations for the binary choice model are given in Cosslett (1987), and the efficient score function $\tilde{\ell}_{\alpha_0, F_0}$ is equal to

$$\tilde{\ell}_{\alpha_0, F_0}(y, \mathbf{x}) = \{x - \mathbb{E}(X | \boldsymbol{\alpha}^T \mathbf{x} = \boldsymbol{\alpha}^T \mathbf{x})\} f(\boldsymbol{\alpha}^T \mathbf{x}) \left\{ \frac{y}{F(\boldsymbol{\alpha}^T \mathbf{x})} - \frac{1-y}{1-F(\boldsymbol{\alpha}^T \mathbf{x})} \right\}, \quad (5)$$

where $f = F'$. For the current status model, the efficient score resembles the efficient score for the binary choice model except that Y and $\boldsymbol{\alpha}^T \mathbf{x}$ are replaced by Δ and $t - \boldsymbol{\alpha}^T \mathbf{x}$, respectively (see, e.g., Huang & Wellner, 1993 or Murphy, van der Vaart, & Wellner, 1999).

The \sqrt{n} -consistent estimators with an asymptotic normal distribution and n times the limiting variance equal to the inverse of

$$\mathbb{E}(\tilde{\ell}_{\alpha, \psi}(\mathbf{X}, Y) \tilde{\ell}_{\alpha, \psi}(\mathbf{X}, Y)^T)$$

are called efficient estimators of α_0 . In the single-index model without the monotonicity constraint on ψ_0 , the efficient estimators of α_0 have been constructed in Ichimura (1993) and Delecroix et al. (2003). Klein and Spady (1993) developed an efficient quasi-maximum likelihood estimator for the binary choice model. An efficient estimate for the current status linear regression model based on a penalized maximum likelihood procedure is proposed in Murphy et al. (1999). Inspired by this penalized estimate, an efficient penalized least squares estimate (PLSE) is constructed for the single-index model with convex link function ψ_0 in Kuchibhotla and Patra (2017).

1.3 | Single-index regression parameter estimators

Several estimators have been proposed in the literature that can be classified into different groups based on the estimation algorithm. Most estimators require a nonparametric estimator for ψ_0 . Often, smoothing procedures, such as kernel smoothers or spline functions, are used to avoid discontinuous criterion functions. An example of this type is the (weighted)-semiparametric least squares estimator (SLSE), which corresponds to minimizing the sum of squares $\sum_{i=1}^n (Y_i - \psi_0(\boldsymbol{\alpha}_0^T \mathbf{X}_i))^2$, when ψ_0 is estimated using a kernel estimator (depending on some bandwidth h). Härdle et al. (1993) extended the SLSE by minimizing the sum of squares over $(\boldsymbol{\alpha}, h)$ simultaneously to obtain optimal smoothing. Instead of estimating ψ_0 by a kernel smoother, spline smoothing is considered in Yu and Ruppert (2002) and Antoniadis, Grégoire, and McKeague (2004), among others. The average derivative estimator proposed by Hristache et al. (2001) results in a direct estimation of the regression parameter α_0 and, therefore, avoids solving a hard optimization problem for M -estimators. The idea of the average derivative method is to estimate the gradient $\boldsymbol{\alpha}_0^T \psi_0'(\boldsymbol{\alpha}_0^T \mathbf{x})$ of the link function using local linear smoothing techniques. A similar approach is considered for the minimum average variance estimator proposed in Xia and Härdle (2006). Smoothing techniques are needed to allow for an efficient estimation of α_0 in the single-index model.

Examples of M -estimators that are not based on an estimate of ψ_0 are Manski's maximum rank estimator (MRE) (Manski, 1975) for the binary choice model, the maximum rank correlation estimator (MRCE) proposed by Han (1987), and the rank estimators proposed by Cavanagh and Sherman (1998) for a more general generalized regression model under monotonicity constraints. The convergence rate for Manski's estimator is in contrast to the other estimators discussed in this section, somewhat disappointingly $n^{1/3}$ instead of the usual \sqrt{n} -rate.

1.4 | Aim of the paper

In this paper, we discuss the behavior of regression parameter estimators in the monotone single-index model. All the estimators are obtained using tuning-parameter-free algorithms and are derived from non-smooth and non-convex criterion functions. In Section 2, we distinguish between two different classes of estimators.

- (a) Estimators with an unknown limiting distribution that depend on the behavior of the piecewise constant, monotone least squares estimator (LSE) of the link function.
- (b) Estimators with a known limiting distribution that converge at the parametric rate to the true regression parameters.

Within the first class of estimators, we first discuss the profiled LSE of the regression parameter in Section 2.1. It is proved in Balabdaoui, Durot, and Jankowski (2016) that the LSE converges at least at the cube-root n rate, but its limiting distribution is still an open problem. Inspired by the rank estimator proposed in Aragón and Quiroz (1995) for the current status model, we also propose a new estimator in this class in Section 2.5.

In Section 2.2, we explain how the \sqrt{n} -consistent estimators proposed in Balabdaoui, Groeneboom, and Hendrickx (2017) for the monotone single-index model can be derived from a score approach. For the second class of estimators, we also describe the rank estimators proposed by Han (1987) and Cavanagh and Sherman (1998) in Section 2.3 and Section 2.4, respectively.

The remainder of Section 2 is devoted to the following.

- The asymptotic properties of the estimators. In particular, we describe a general approach in Section 2.6 for proving the asymptotic normal distribution of the estimators of the second class (b), and we discuss, in Section 2.7, the difficulties that arise when one wants to derive the limiting distribution of the estimators of the first class (a).
- The computation of the estimators.

The quality of the estimators is illustrated via simulations and a real data example in Section 3 and Section 4, respectively. Although smoothing is necessary to obtain efficient estimators in the single-index model, we want to point out that smoothing should not be the main concern when interest is in estimating the finite-dimensional regression parameter. A slight loss of efficiency is an acceptable price to pay for a tuning-free parameter procedure that is computationally more attractive than efficient procedures (since no smoothing parameter selection is needed). Efficient estimators are moreover often based on smoothness conditions that are stronger than the conditions needed when smoothing techniques are avoided.

2 | ESTIMATORS OBTAINED WITHOUT SMOOTHING TECHNIQUES

In this section, we describe different estimators for α_0 based on a random sample $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ of size n from (\mathbf{X}, Y) , where $\mathbb{E}(Y|\mathbf{X}) = \psi_0(\alpha_0^T \mathbf{X})$ and ψ_0 belongs to the class \mathcal{M} of monotone functions on \mathbb{R} . To ensure identifiability of the single-index model, we assume that α_0 is a vector of regression parameters belonging to the $(d - 1)$ -dimensional sphere $\mathcal{S} := \{\alpha \in \mathbb{R}^d : \|\alpha\| = 1\}$.

To illustrate the criterion functions associated with the different estimators, we consider a simulated data sample from the model

$$Y = \exp\left(X_1/\sqrt{2} + X_2/\sqrt{2}\right) + \varepsilon, \quad X_1, X_2 \sim U[-1, 1] \text{ and } \varepsilon \sim N(0, 1). \quad (6)$$

For each of the estimators, we include figures for the criterion function as a function of θ , where θ is defined by $(\alpha_1, \alpha_2) = (\cos(\theta), \sin(\theta))$. In model (6), the true parameter value is $\theta_0 = \pi/4$.

2.1 | The least squares estimator

Consider the sum of squared errors

$$S_n(\alpha, \psi) = \frac{1}{n} \sum_{i=1}^n \{Y_i - \psi(\alpha^T \mathbf{X}_i)\}^2, \quad (7)$$

which can be computed for any pair $(\alpha, \psi) \in S_{d-1} \times \mathcal{M}$. The LSE $(\hat{\alpha}_n, \hat{\psi}_n)$ is defined by

$$(\hat{\alpha}_n, \hat{\psi}_n) := \arg \min_{\alpha \in S, \psi \in \mathcal{M}} S_n(\alpha, \psi). \quad (8)$$

The LSE can be obtained as follows. For a fixed $\alpha \in S_{d-1}$, order the values $\alpha^T \mathbf{X}_1, \dots, \alpha^T \mathbf{X}_n$ in increasing order and arrange Y_1, \dots, Y_n accordingly. As ties are not excluded, let $m = m_\alpha$ be the number of distinct projections among $\alpha^T \mathbf{X}_i$ and $Z_1^\alpha < \dots < Z_m^\alpha$ be the corresponding ordered values. For $i = 1, \dots, m$, let

$$n_i^\alpha = \sum_{j=1}^n 1_{\{\alpha^T \mathbf{X}_j = Z_i^\alpha\}} \quad \text{and} \quad Y_i^\alpha = \sum_{j=1}^n Y_j 1_{\{\alpha^T \mathbf{X}_j = Z_i^\alpha\}} / n_i^\alpha.$$

Then, well-known results from the isotonic regression theory imply that the functional $\psi \mapsto S_n(\alpha, \psi)$ is minimized by the left derivative of the greatest convex minorant of the cumulative sum diagram

$$\left\{ (0, 0), \left(\sum_{j=1}^i n_j^\alpha, \sum_{j=1}^i n_j^\alpha Y_j^\alpha \right), i = 1, \dots, m \right\}. \quad (9)$$

See, for example, theorem 1.1 in Barlow, Bartholomew, Bremner, and Brunk (1972) or theorem 1.2.1 in Robertson, Wright, and Dykstra (1988). By strict convexity of $\psi \mapsto S_n(\psi, \alpha)$, the minimizer is unique at the distinct projections. We define the nonparametric least squares estimator (LSE(ψ)) $\hat{\psi}_{n\alpha}$ by the monotone function that takes the values of this minimizer at the distinct projections and is extended to a right-continuous step function outside the set of those projections. Figure 1 shows a picture of the LSE of $\hat{\psi}_{n\alpha_0}$ for the model described in (6). The LSE $(\hat{\alpha}_n, \hat{\psi}_n) := (\hat{\alpha}_n, \hat{\psi}_{n\hat{\alpha}_n})$ is next obtained by maximizing the map $\alpha \mapsto S_n(\alpha, \hat{\psi}_{n\alpha})$ over all $\alpha \in S$.

Since the cumulative sum diagram in (9) only changes when the ordering in $\alpha^T \mathbf{X}_i$ changes, the vector $(\hat{\psi}_{n\alpha}(\alpha^T \mathbf{X}_1), \dots, \hat{\psi}_{n\alpha}(\alpha^T \mathbf{X}_n))^T$ will be the same for all α for which the ranks of the $\alpha^T \mathbf{X}_i$ are the same. As a consequence, the criterion function $\alpha \mapsto S_n(\alpha, \hat{\psi}_{n\alpha})$ is piecewise constant, and the LSE $(\hat{\alpha}_n, \hat{\psi}_{n\hat{\alpha}_n})$ is not unique.

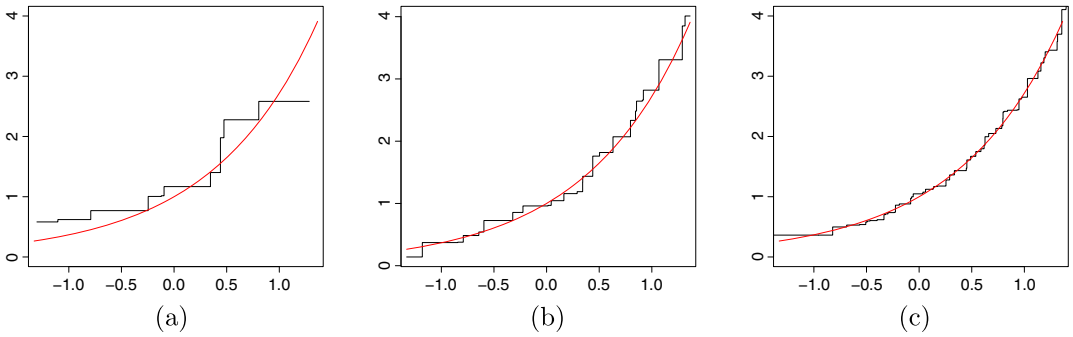


FIGURE 1 The LSE $\hat{\psi}_{n,\alpha_0}$ (black, step-wise) and the true function $\psi_0(x) = \exp(x)$ (red, solid) in model (6) for a sample of sizes (a) $n = 100$, (b) $n = 1000$, and (c) $n = 10000$

2.1.1 | The maximum likelihood estimator for the binary choice model and the current status linear regression model

Instead of considering the LSE for the binary choice model or the current status model, we consider maximizing the log likelihood of the data. In particular, for the binary choice model, the MLE $(\hat{\alpha}_n, \hat{F}_n)$ is defined by

$$(\hat{\alpha}_n, \hat{F}_n) := \arg \max_{\alpha \in \mathcal{S}, F \in \mathcal{F}} L_n(\alpha, F),$$

where \mathcal{F} is the set of all distribution functions on \mathbb{R} and

$$L_n(\alpha, F) := \frac{1}{n} \sum_{i=1}^n [Y_i \log F(\alpha^T \mathbf{X}_i) + (1 - Y_i) \log \{1 - F(\alpha^T \mathbf{X}_i)\}]. \tag{10}$$

Using profiled log likelihood, we first obtain the minimizer of the map $F \mapsto L_n(\alpha, F)$ over \mathcal{F} , which is, in fact, the same function as the minimizer of $S_n(\alpha, F)$. For the current status model, we replace Y_i and $\alpha^T \mathbf{X}_i$ in the expression for L_n given in (10) by Δ_i and $T_i - \alpha^T \mathbf{X}_i$ and maximize over $\mathbb{R} \times \mathcal{F}$ instead of $\mathcal{S} \times \mathcal{F}$ due to the identifiability of the current status regression model.

Note that in a homoscedastic model with normal error terms, the MLE for the single-index model (1) is equivalent to the LSE. The asymptotic properties of the LSE in (8) for the monotone single-index model and the MLE in (10) for the binary choice model and the current status model are discussed in Balabdaoui et al. (2016), Cosslett (1983), and Murphy et al. (1999), respectively. All estimators converge at rate $n^{1/3}$, but their limiting distribution is still an open problem.

2.2 | The simple score estimator

Balabdaoui et al. (2017) developed a Z -estimator based on the derivative of the sum of squares $S_n(\alpha, \hat{\psi}_{n\alpha})$, ignoring the non-differentiability of the LSE $\hat{\psi}_{n\alpha}$. They first consider a local parametrization \mathbb{S} mapping a subset of \mathbb{R}^{d-1} to the sphere \mathcal{S} . Examples of such parametrizations are the map

$$(\beta_1, \beta_2, \dots, \beta_{d-1}) \mapsto \left(\sqrt{1 - \beta_2^2 - \dots - \beta_{d-1}^2}, \beta_2, \dots, \beta_d \right)^T$$

or the spherical coordinate system

$$(\beta_1, \beta_2, \dots, \beta_{d-1}) \mapsto (\cos(\beta_1), \sin(\beta_1) \cos(\beta_2), \sin(\beta_1) \sin(\beta_2) \cos(\beta_3), \dots, \sin(\beta_1) \cdots \sin(\beta_{d-2}) \cos(\beta_{d-1}), \sin(\beta_1) \cdots \sin(\beta_{d-2}) \sin(\beta_{d-1}))^T.$$

The simple score estimator (SSE) is next obtained by a zero-crossing (in β) of

$$Z_n(\beta) := \frac{1}{n} \sum_{i=1}^n (J_{\mathbb{S}}(\beta))^T \mathbf{X}_i \{Y_i - \hat{\psi}_{n\alpha}(\alpha^T \mathbf{X}_i)\}, \quad (11)$$

where $\alpha = \mathbb{S}(\beta)$ and $J_{\mathbb{S}} \in \mathbb{R}^{d \times d-1}$ is the matrix of the partial derivatives of $\mathbb{S} : \mathbb{R}^{d-1} \mapsto \mathcal{S} : \beta \mapsto \mathbb{S}(\beta) = \alpha$. Since the LSE $\hat{\psi}_{n\alpha}$ is the same for different α values, the score criterion (11) will have discontinuities and an exact root of (11) does not always exist. The estimator is therefore defined by $\hat{\alpha}_n = \mathbb{S}(\hat{\beta}_n)$, where $\hat{\beta}_n$ is a point in \mathbb{R}^{d-1} such that each component of Z_n crosses through zero at $\hat{\beta}_n$.

The SSE is, in fact, a generalization of the related estimator for the current status linear regression model proposed in Groeneboom and Hendrickx (2017) and defined by a zero-crossing of

$$Z'_n(\alpha) := \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \{\Delta_i - \hat{F}_{n\alpha}(T_i - \alpha^T \mathbf{X}_i)\}, \quad (12)$$

where $\hat{F}_{n\alpha}$ is the MLE for F_0 when α is fixed.

2.3 | The maximum rank correlation estimator

Han's MRCE is motivated by the fact that $Y_i \geq Y_j$ is more likely than $Y_i < Y_j$ when $\alpha_0^T \mathbf{X}_i \geq \alpha_0^T \mathbf{X}_j$ if ψ_0 is increasing. The MRCE is defined by the maximizer of

$$H_n(\alpha) := \frac{1}{n(n-1)} \sum_{i \neq j} \{Y_i > Y_j\} \{\alpha^T \mathbf{X}_i > \alpha^T \mathbf{X}_j\}. \quad (13)$$

In contrast to the LSE and the SSE, estimation of the unknown link function ψ_0 is not considered with the MRCE.

2.4 | The maximum rank estimator

Inspired by the MRCE, Cavanagh and Sherman (1998) developed a new class of rank estimators defined by the maximizer of

$$R_n(\alpha) := \frac{1}{n(n-1)} \sum_{i \neq j} M(Y_i) \{\alpha^T \mathbf{X}_i > \alpha^T \mathbf{X}_j\}, \quad (14)$$

where M denotes an increasing function on \mathbb{R} . In this paper, we investigate the behavior of the estimator when M is equal to the identity function, that is, $M(y) = y$, and refer to this estimator as the MRE. Since the responses in the binary choice model and the current status model are binary, it can be shown (see Appendix, Section A.1) that the MRCE and the MRE are equivalent in these models. The behavior of the map $\alpha \mapsto H_n(\alpha)$ and the map $\alpha \mapsto R_n(\alpha)$ are similar, and we do not include pictures for the latter mapping.

2.5 | The maximum rank estimator using the least squares estimator of ψ_0

Aragón and Quiroz (1995) proposed two regression parameter estimators for the current status linear regression model based on the ranks of the observations $T_i - \alpha^T \mathbf{X}_i$. The first estimator coincides with the MRE. The second estimator is defined by the maximizer of

$$\sum_{i \neq j} \hat{F}_{n\alpha}(T_i - \alpha^T \mathbf{X}_i) \{T_i - \alpha^T \mathbf{X}_i > T_j - \alpha^T \mathbf{X}_j\},$$

TABLE 1 Asymptotic variance. Monotone single-index model

Method	Σ
SSE	$\mathbb{E} \left[\{Y - \psi_0(\alpha_0^T \mathbf{X})\}^2 \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T \right]$
MRCE	$\mathbb{E} \left[\{2F_0(Y - \psi_0(\alpha^T \mathbf{X})) - 1\}^2 \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T g_0(\alpha_0^T \mathbf{X})^2 \right]$
MRE	$\mathbb{E} \left[\{Y - \psi_0(\alpha^T \mathbf{X})\}^2 \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T g_0(\alpha_0^T \mathbf{X})^2 \right]$
\mathbf{V}	
SSE	$\mathbb{E} \left[\psi_0'(\alpha_0^T \mathbf{X}) \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T \right]$
MRCE	$\mathbb{E} \left[2\psi_0'(\alpha_0^T \mathbf{X}) f_0(Y - \psi_0(\alpha^T \mathbf{X})) \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T g_0(\alpha_0^T \mathbf{X}) \right]$
MRE	$\mathbb{E} \left[\psi_0'(\alpha_0^T \mathbf{X}) \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T g_0(\alpha_0^T \mathbf{X}) \right]$

Note. SSE = simple score estimator; MRCE = maximum rank correlation estimator; MRE = maximum rank estimator.

where $\hat{F}_{n\alpha}$ is the MLE (see (10)) for fixed α . This motivates us to investigate the behavior of the estimator for the monotone single-index model, referred to as the LS-MRE, defined by the maximizer of

$$A_n(\alpha) := \frac{1}{n(n-1)} \sum_{i \neq j} \hat{\psi}_{n\alpha}(\alpha^T \mathbf{X}_i) \{ \alpha^T \mathbf{X}_i > \alpha^T \mathbf{X}_j \}, \tag{15}$$

where $\hat{\psi}_{n\alpha}$ is the LSE for fixed α . To the best of our knowledge, this estimator has not been studied before and the asymptotic limiting distribution is still unknown. Since the LS-MRE is similar to the LSE, an M -estimator that involves the nonparametric LSE for ψ_0 , it can be expected that similar theoretical issues appear when deriving the limiting behavior for both estimators.

2.6 | Asymptotic behavior

It has been shown in Balabdaoui et al. (2017) for the SSE, in Sherman (1993) for the MRCE, and in Cavanagh and Sherman (1998) for the MRE that these estimators are \sqrt{n} -consistent and have an asymptotic normal distribution with an asymptotic variance that is larger than the efficient variance. As pointed out in a footnote on p. 361 of Cavanagh and Sherman, the expression for the asymptotic variance of the MRCE given in theorem 4 is only correct up to a factor 4. Unfortunately, the same mistake for the MRE was made in the expression for the asymptotic variance of the MRE given in theorem 2 of Cavanagh and Sherman.

Although no proof for the MRCE and the MRE has been published, we can prove that

$$\sqrt{n}(\hat{\alpha}_n - \alpha_0) \rightarrow_d N(\mathbf{0}, \mathbf{V}^- \Sigma, \mathbf{V}^-), \tag{16}$$

where \mathbf{V}^- is the Moore–Penrose inverse of \mathbf{V} ; a sketch of the proof of (16) is given in the Appendix, Section A.2. The reason that we have to consider generalized inverses is that the normal limiting distributions are concentrated on the $(d - 1)$ -dimensional subspace, orthogonal to α_0 , and therefore degenerate. This is also clear from its covariance matrix $\mathbf{V}^- \Sigma \mathbf{V}^-$, which is a matrix of rank $d - 1$. The expressions for \mathbf{V} and Σ are summarized in Table 1 for the monotone single-index model and in Table 2 for the current status linear regression model.

The limiting distributions of the LSE and the LS-MRE are still unknown. Figures 2 and 3 show a more irregular behavior of the criterion functions for the LSE and the LS-MRE compared to the smoother criterion functions for the SSE and the MRCE, shown in Figures 4 and 5.

TABLE 2 Asymptotic variance. Current status linear regression model

Method	Σ
SSE	$\mathbb{E} \left[F_0(T - \alpha_0^T \mathbf{X}) \{1 - F_0(T - \alpha_0^T \mathbf{X})\} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} T - \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} T - \alpha_0^T \mathbf{X}) \}^T \right]$
MR(C)E	$\mathbb{E} \left[F_0(T - \alpha_0^T \mathbf{X}) \{1 - F_0(T - \alpha_0^T \mathbf{X})\} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} T - \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} T - \alpha_0^T \mathbf{X}) \}^T g_0(T - \alpha_0^T \mathbf{X})^2 \right]$
V	
SSE	$\mathbb{E} \left[f_0(T - \alpha_0^T \mathbf{X}) \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T \right]$
MR(C)E	$\mathbb{E} \left[f_0(T - \alpha_0^T \mathbf{X}) \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \} \{ \mathbf{X} - \mathbb{E}(\mathbf{X} \alpha_0^T \mathbf{X}) \}^T g_0(T - \alpha_0^T \mathbf{X}) \right]$

Note. SSE = simple score estimator; MRCE = maximum rank correlation estimator; MRE = maximum rank estimator.

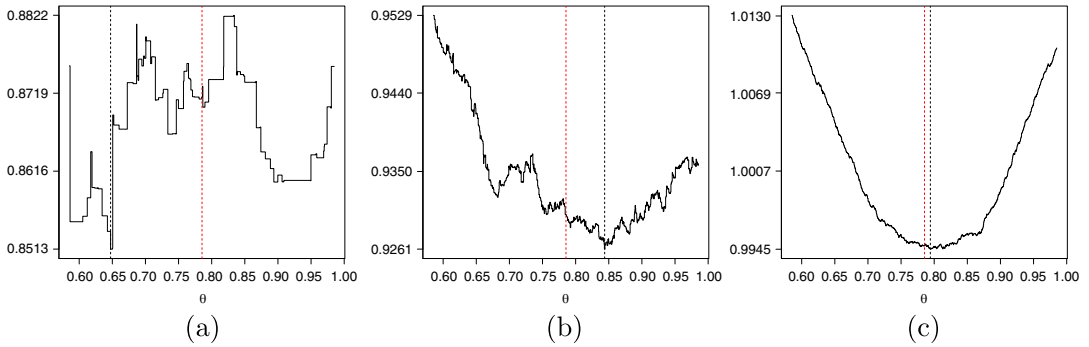


FIGURE 2 The map $\theta \mapsto S_n((\cos(\theta), \sin(\theta))^T, \hat{\psi}_{n,\alpha})$ (black, solid) in model (6) for a sample of sizes (a) $n = 100$, (b) $n = 1000$, and (c) $n = 10000$. The vertical reference lines indicate the position of the minimizer (black, dotted) and the true parameter value $\theta_0 = \pi/4$ (red, dotted)

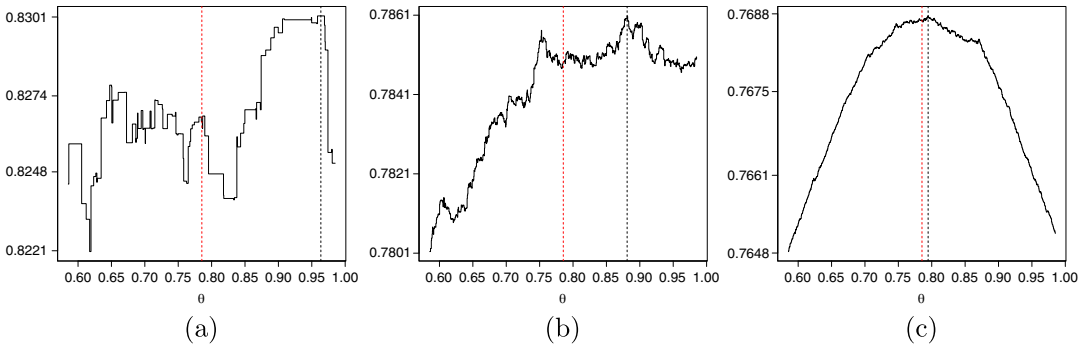


FIGURE 3 The map $\theta \mapsto A_n((\cos(\theta), \sin(\theta))^T)$ (black, solid) in model (6) for a sample of sizes (a) $n = 100$, (b) $n = 1000$, and (c) $n = 10000$. The vertical reference lines indicate the position of the maximizer (black, dotted) and the true parameter value $\theta_0 = \pi/4$ (red, dotted)

2.7 | Difficulties with the LSE and the LS-MRE

Deriving the limiting distributions for the LSE and the LS-MRE is challenging. One of the difficulties arises from the non-differentiability of the LSE $\hat{\psi}_{n\hat{\alpha}}$ for ψ_0 appearing in the criterion functions S_n and A_n . This is, for example, not the case with the efficient SLSE proposed in Ichimura (1993), where the criterion function is given by S_n defined in (7) but with $\hat{\psi}_{n\hat{\alpha}}$ replaced by a kernel estimate

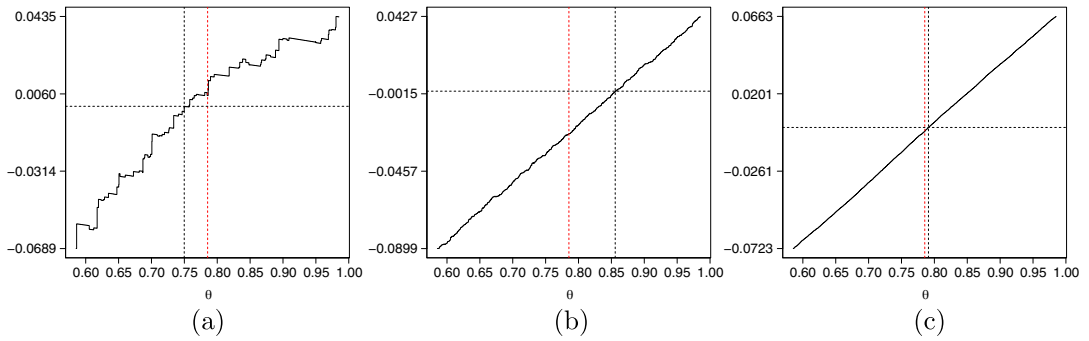


FIGURE 4 The map $\theta \mapsto Z_n((\cos(\theta), \sin(\theta)^T))$ (black, solid) in model (6) for a sample of sizes (a) $n = 100$, (b) $n = 1000$, and (c) $n = 10000$. The vertical reference lines indicate the position of the zero-crossing (black, dotted) and the true parameter value $\theta_0 = \pi/4$ (red, dotted)

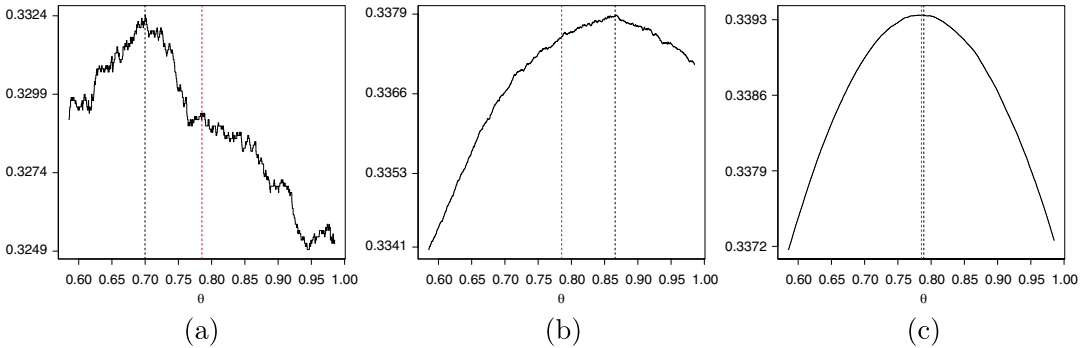


FIGURE 5 The map $\theta \mapsto H_n((\cos(\theta), \sin(\theta)^T))$ (black, solid) in model (6) for a sample of sizes (a) $n = 100$, (b) $n = 1000$, and (c) $n = 10000$. The vertical reference lines indicate the position of the maximizer (black, dotted) and the true parameter value $\theta_0 = \pi/4$ (red, dotted)

that is two times continuously differentiable with respect to α . By considering a Z-estimator instead of an M-estimator, this non-differentiability is somehow circumvented with the SSE.

As discussed in Groeneboom and Hendrickx (2017), the “canonical” approach to proofs that argmax or argmin estimates of α_0 are \sqrt{n} -consistent has been provided by Sherman (1993) for the models, considered in his paper. His theorem 1 says that $\|\hat{\alpha}_n - \alpha_0\| = O_p(n^{-1/2})$, where $\|\cdot\|$ denotes the Euclidean norm, if $\hat{\alpha}_n$ is the maximizer of a criterion function $\Gamma_n(\alpha)$, with population equivalent $\Gamma(\alpha)$ and

(a) there exists a neighborhood N of α_0 and a constant $k > 0$ such that

$$\Gamma(\alpha) - \Gamma(\alpha_0) \leq -k\|\alpha - \alpha_0\|^2,$$

for $\alpha \in N$, and

(b) uniformly over $o_p(1)$ neighborhoods of α_0 ,

$$\Gamma_n(\alpha) - \Gamma_n(\alpha_0) = \Gamma(\alpha) - \Gamma(\alpha_0) + O_p\left(\|\alpha - \alpha_0\|/\sqrt{n}\right) + o_p\left(\|\alpha - \alpha_0\|^2\right) + O_p(n^{-1}).$$

Moreover, assuming $\alpha_0 = \mathbf{0}$ and $\Gamma(\alpha_0) = \mathbf{0}$, he continues by making the assumption that

$$\Gamma_n(\alpha) = -\frac{1}{2}\alpha^T \mathbf{V}\alpha + n^{-1/2}\alpha^T \mathbf{W}_n + o_p(n^{-1}),$$

where \mathbf{V} is a positive definite matrix and \mathbf{W}_n converges in distribution to a normal distribution. Under these circumstances, $\sqrt{n}\alpha_n$ also converges to a normal distribution.

If we try to apply this to our situation, we first have to deal with the fact that we have to use a parametrization of the type introduced in Section 2.2 in order to consider full neighborhoods of β_0 and nondegenerate matrices in \mathbb{R}^{d-1} in the single-index model. But then, for the profile LSE estimator $\hat{\beta}_n$, it is not clear that an expansion as given in Sherman (1993) will hold. We seem to get inevitably an extra term of order $O_p(n^{-2/3})$ in (b) (with α_n and α_0 replaced by β_n and β_0), which does not fit into this framework.

On the other hand, in the expansion of our score function (11), we get that this function is in first order the sum of a term of the form

$$\phi'(\beta_0)(\beta - \beta_0),$$

where ϕ' is the matrix, representing the total derivative of the population equivalent score function, and a term \mathbf{W}_n of order $O_p(n^{-1/2})$, which gives

$$\hat{\beta}_n - \beta_0 \sim -\phi'(\beta_0)^{-1}\mathbf{W}_n = O_p(n^{-1/2}),$$

and, here, extra terms of order $O_p(n^{-2/3})$ do not hurt.

Li and Zhang (1998), in their paper on smooth U-statistics estimators for the regression parameter in the current status linear regression model, conjecture in their remark 2.4 that the profile MLE (which is similar to the LSE) will be \sqrt{n} -convergent, but not efficient. We could not follow their argument for this conjecture. Aragón and Quiroz (1995) have a similar conjecture for their second estimate, based on the nonparametric MLE, but again, the argument for this conjecture in their comment (i) in section 4 is not at all convincing. In fact, our simulations show that the simple argmin or argmax estimates seem to be always inferior to the simple score estimates, and the pictures of the argmin-type estimates as a function of the parameter show that the behavior is a lot more irregular than the behavior of the score estimates. Hence, even if the argmin-type or argmax-type estimators would be \sqrt{n} -convergent, the score estimators have a better behavior.

2.8 | Computation

None of the discontinuous criterion functions in Sections 2.1–2.5 are convex. This makes the computation of the estimators difficult. Standard optimization methods for convex loss functions cannot be used. The discreteness of the criterion functions moreover excludes methods that take derivative information into account since this derivative is often not defined.

For the computation of the M -estimators (LSE, MRCE, MRE, and LS-MRE), we wrote C++ programs, using the pattern search algorithm proposed by Hooke and Jeeves (1961). The latter optimization method searches for the minimum of a loss function without requiring a gradient. Based on an initial starting value, $2d$ local searches of a certain step size are made in each direction to see if a lower loss function is obtained. The first move is made in the direction of the previous move. If no function decrease is found, the step size is reduced. The procedure is iterated until a convergence criterion is satisfied. The convergence of this pattern search algorithm is discussed in Torczon (1997). The root of the set of $d - 1$ score functions in Section 2.2 coincides with the minimizer of the sum of squared component score functions so that this minimization approach

can also be used to obtain the SSE. We recognize the importance of good starting values and improvements of the current approach are worth studying in future computational research.

3 | SIMULATIONS

To evaluate the finite sample behavior of the different estimators introduced in Section 2, we simulate $N = 5000$ data sets from the model

$$Y = \psi_0(\boldsymbol{\alpha}_0^T \mathbf{X}) + \varepsilon, \quad (17)$$

where $\psi_0(x) = x + x^3$, $\alpha_{0i} = 1/\sqrt{3}$, $i = 1, 2, 3$, and $\varepsilon \sim N(0, 1)$, independent of \mathbf{X} . We consider two different distributions for the covariate vector \mathbf{X} , $X_i \stackrel{i.i.d.}{\sim} U[0, 1]$ and $X_i \stackrel{i.i.d.}{\sim} N(0, 1)$ for $i = 1, 2, 3$.

Tables 3 and 4 show the mean and n times the covariance matrix of the estimates for sample sizes $n = 100, 500, 1000, 5000$, and 100000 for the uniform, respectively normal, simulation setting. For the SSE, the MRCE, and the MRE, we calculated the asymptotic variances given in Table 1 to which n times the covariance matrix should converge. We however note that only the uniform model satisfies the assumptions needed to prove (16). The last column in Tables 3 and 4 contains the distance between n times the covariance matrix of the estimates and the matrix $\mathbf{V}^{-1}\boldsymbol{\Sigma}\mathbf{V}^{-1}$ obtained by summing up the squared distance of the corresponding matrix elements. The results for n times the variance of the estimates of α_3 are visualized in Figure 6.

For both simulation settings, the results show the convergence of n times the covariance matrix towards the asymptotic values for the SSE, MRCE, and MRE. The convergence rate is faster for the SSE than for the MRCE and the MRE. The asymptotic values are smallest for the SSE in these models, with only a small difference for the uniform setting but a larger difference in the normal setting where the asymptotic values of the MRCE and the MRE are substantially larger than the ones for the SSE.

For the LS-MRE, n times the covariance matrices increase with increasing sample size, suggesting a slower convergence rate than the parametric \sqrt{n} -rate for this estimator. Table 3 also shows a similar increase for the LSE in the uniform model, whereas a decrease of n times the covariance matrix for the LSE is shown in Table 4 for the normal setting. The LSE even performs better than the MRCE and the MRE in the latter simulation model.

Finally, we also compared the estimates in the uniform model with the PLSE proposed by Kuchibhotla and Patra (2017) and the EFM estimate proposed by Cui, Härdle, and Zhu (2011). The function `simestgcv` available in the R-packages `simest` was used to obtain the PLSE.

The algorithm described on p. 1670 of Cui et al. (2011) was implemented by us in a C++ program, with a fixed tuning parameter M for the EDF approach. The computation time is considerably longer than the time required for the methods discussed in Section 2 (this observation also holds for the PLSE estimate proposed in Kuchibhotla & Patra, 2017). Moreover, the fixed-point algorithm, used in the second step of the algorithm of Cui et al. (2011), has an oscillating behavior and will certainly not converge in any monotone way; one has to wait whether it will enter at a certain point of the iterations into a sufficiently small ball around a fixed point. The version of the algorithm, kindly sent to us by Dr. Cui, depends crucially on the R algorithm `smooth.lf` in the R package `locfit`, and it is not so easy to see how the first step of the algorithm is handled by this package. We used our own C++ implementation of the algorithm, but noticed that the R implementation had a similar behavior. The R implementation of Cui et al. (2011) also uses a fixed tuning parameter M for the EDF approach. In both implementations, one has to discard non-converging runs.

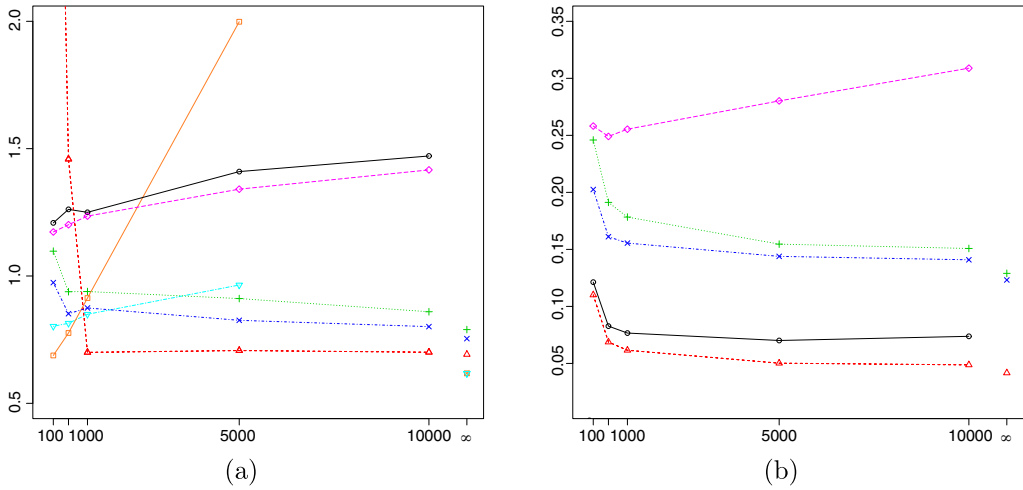


FIGURE 6 n times the variance of α_3 as a function of the sample size n for the simulation model (a) with $X_i \sim U[0, 1]$ and (b) with $X_i \sim N(0, 1)$ for the LSE (solid, black, \circ), SSE (dashed, red, Δ), MRCE (dotted, green, $+$), MRE (dashed-dotted, blue, \times), LS-MRE (long-dashed, pink, \diamond), EDF (two-dashed, light blue, ∇), and PLSE (solid, orange, \square). The point at ∞ represents the asymptotic values

Therefore, we do not report results for the sample size $n = 10000$ and simulated only $N = 2500$ data sets for the PLSE with $n = 5000$. Boxplots of $\sum_{j=0}^3 (\hat{\alpha}_j - \alpha_{0j})^2 / 3$, shown in Figure 7, illustrate that the PLSE and EFM estimates perform better than the SSE, MRCE, and MRE for smaller sample sizes. As the sample size increases, the results for the efficient but computational intensive methods are no longer superior, and the best performance is obtained with the SSE. The results for the PLSE and EDF estimates depend furthermore on smoothing parameters that need to be selected carefully. Figure 6 clearly shows that n times the variance increases for the PLSE with increasing sample size, in contrast to the efficient convergence rate. This illustrates again that, in practice, methods involving smoothing techniques are not necessarily a better choice than \sqrt{n} -consistent parameter-free methods, especially for larger sample sizes where the computation cost is enormous.

We conclude that it is worthwhile to consider parameter-free methods for estimation in the monotone single-index model. The additional complexity (due to the smoothing parameter) does not necessarily result in better performances for efficient estimates. The increased computation time is only worthwhile when the sample size is small. The SSE is preferred for larger samples and moreover achieves better performances than the rank estimators (MRCE and MRE). The experiments in the normal model were in favor of the parametric \sqrt{n} -rate for the LSE, whereas the uniform trials suggested a slower convergence rate. Even if the LSE leads at all to an \sqrt{n} -consistent estimate, its performance remains inferior to the score procedure in Section 2.2. Nevertheless, it remains an interesting topic to understand the behavior of the LSE in the monotone single-index model.

4 | REAL DATA EXAMPLE

In this section, we apply the estimation techniques to the ozone data (Chambers, Cleveland, Kleiner, & Tukey, 1983). The data set contains observations on the ozone concentration for

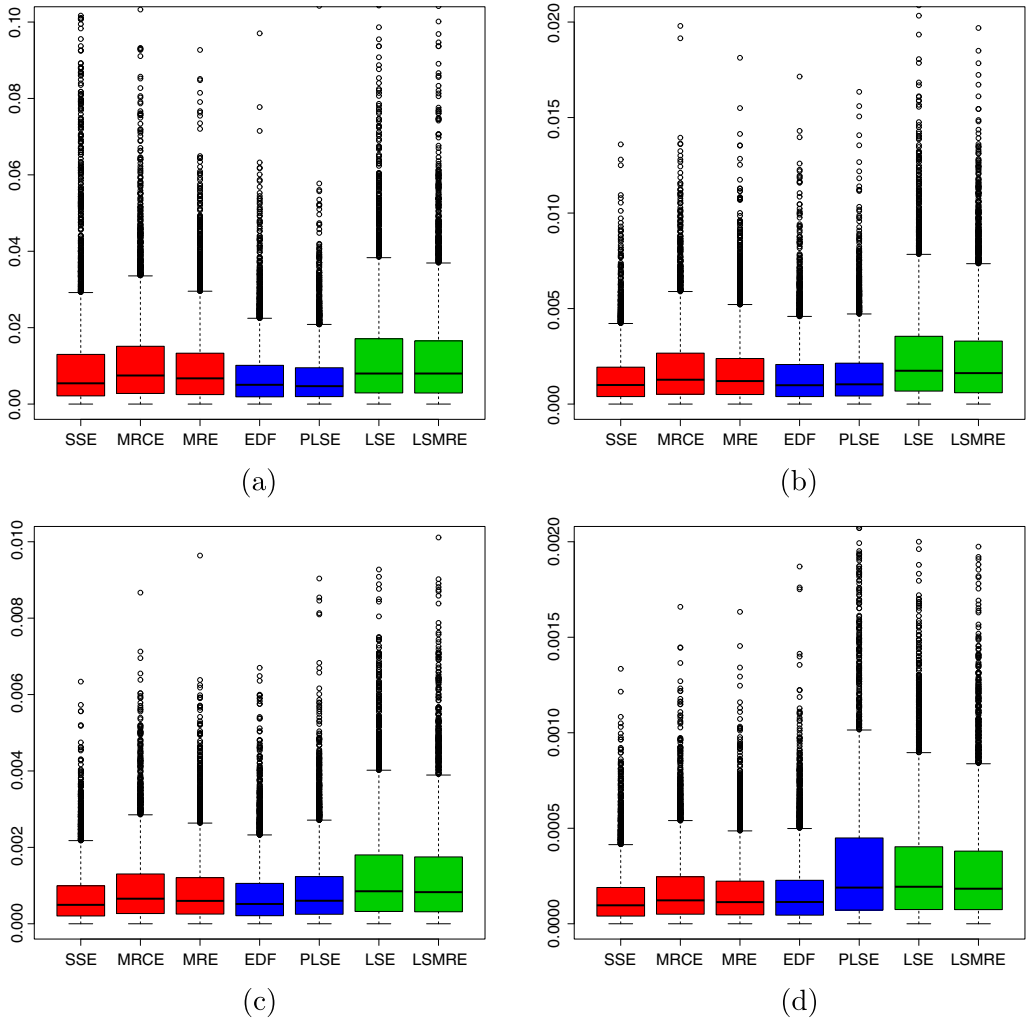


FIGURE 7 Boxplots of $\sum_{j=0}^3 (\hat{\alpha}_j - \alpha_{0j})^2 / 3$ for the model with $X_i \sim U[0, 1]$ for sample sizes (a) $n = 100$, (b) $n = 500$, (c) $n = 1000$, and (d) $n = 5000$. Red boxes correspond to \sqrt{n} -consistent but inefficient methods (SSE, MRCE, and MRE), blue boxes correspond to \sqrt{n} -consistent and efficient methods (EDF and PLSE), and green boxes correspond to methods with an unknown limiting distribution (LSE and LS-MRE)

153 consecutive days between May 1, 1973 and September 30, 1973. We study the relationship between the ozone concentration (Y , ppb) and the meteorological variables, namely, solar radiation (R , Ly), temperature (T , $^{\circ}F$), and wind speed (W , mph), in a subset of the data consisting of 111 complete observations.

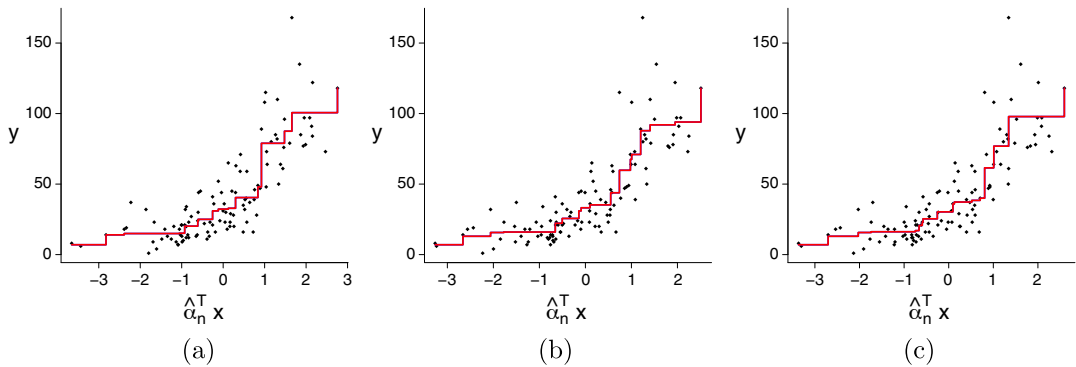
Yu and Ruppert (2002) compared linear, semiparametric, and fully nonparametric fits for the ozone data and concluded that better results were obtained with a single-index model. The data set became a benchmark in the study of single-index models (Alquier & Biau, 2013; Antoniadis et al., 2004; Karabatsos, 2009; and Wang, 2009, among others), and the results, previously presented in the statistical literature, suggest that a monotone fit for the underlying link function is plausible for the ozone data.

For our data analysis, we have scaled the covariates to have mean 0 and variance 1. Table 5 summarizes the results of the regression parameters for the LSE, SSE, MRCE, MRE, and LS-MRE.

TABLE 5 Ozone data

Method	R	T	W
LSE	0.261650	0.673180	-0.691641
SSE	0.288573	0.857762	-0.425406
MRCE	0.371694	0.833361	-0.409088
MRE	0.380572	0.835861	-0.395603
LS-MRE	0.269241	0.828638	-0.490783

Note. LSE = least squares estimator; SSE = simple score estimator; MRCE = maximum rank correlation estimator; MRE = maximum rank estimator.

**FIGURE 8** Ozone data. Scatterplot $(\hat{\alpha}_n^T \mathbf{x}_i, y_i)$ and $\hat{\psi}_{n\hat{\alpha}_n}$ (red, step function) for (a) LSE, (b) SSE, and (c) LS-MRE

The estimate $\hat{\psi}_{n\hat{\alpha}_n}$ of ψ_0 , together with a scatterplot of $(\hat{\alpha}_n^T \mathbf{x}_i, y_i)$, is given in Figure 8 for the LSE, SSE, and LS-MRE. We see that the estimates described in this paper result in similar estimated relationships between the ozone concentration and the meteorological variables.

ACKNOWLEDGEMENTS

The research of the second author was supported by the Research Foundation Flanders (FWO) under Grant 11W7315N. Support from the IAP Research Network P7/06 of the Belgian State (Belgian Science Policy) is gratefully acknowledged. For the simulations, we used the infrastructure of the VSC—Flemish Supercomputer Center, funded by the Hercules Foundation and the Flemish Government—Department EWI.

ORCID

Piet Groeneboom  <http://orcid.org/0000-0001-8027-8114>

Kim Hendrickx  <http://orcid.org/0000-0003-4005-2676>

REFERENCES

- Alquier, P., & Biau, G. (2013). Sparse single-index model. *Journal of Machine Learning Research*, 14, 243–280.
- Antoniadis, A., Grégoire, G., & McKeague, I. W. (2004). Bayesian estimation in single-index models. *Statistica Sinica*, 1147–1164.

- Aragón, J., & Quiroz, A. J. (1995). Rank regression for current status data. *Statistics & Probability Letters*, 24(3), 251–256.
- Balabdaoui, F., Durot, C., & Jankowski, H. (2016). Least squares estimation in the monotone single index model. arXiv preprint arXiv:1610.06026.
- Balabdaoui, F., Groeneboom, P., & Hendrickx, K. (2017). Score estimation in the monotone single index model. arXiv Preprint arXiv:1712.05593.
- Barlow, R., Bartholomew, D., Bremner, J., & Brunk, H. (1972). *Statistical inference under order restrictions. The theory and application of isotonic regression*, Wiley Series in Probability and Mathematical Statistics. London–New York–Sydney: John Wiley & Sons.
- Cavanagh, C., & Sherman, R. P. (1998). Rank estimators for monotonic index models. *Journal of Econometrics*, 84(2), 351–381.
- Chambers, J., Cleveland, W., Kleiner, B., & Tukey, P. (1983). *Graphical methods for data analysis*. Boston: Duxbury Press.
- Cosslett, S. (1987). Efficiency bounds for distribution-free estimators of the binary choice and the censored regression models. *Econometrica*, 55(3), 559–585.
- Cosslett, S. R. (1983). Distribution-free maximum likelihood estimator of the binary choice model. *Econometrica*, 51(3), 765–782.
- Cui, X., Härdle, W. K., & Zhu, L. (2011). The EFM approach for single-index models. *The Annals of Statistics*, 39(3), 1658–1688.
- Delecroix, M., Härdle, W., & Hristache, M. (2003). Efficient estimation in conditional single-index regression. *Journal of Multivariate Analysis*, 86(2), 213–226.
- Groeneboom, P., & Hendrickx, K. (2017). Current status linear regression. *The Annals of Statistics*. Retrieved from <https://arxiv.org/abs/1601.00202>
- Groeneboom, P., & Jongbloed, G. (2014). *Nonparametric estimation under shape constraints*. Cambridge: Cambridge University Press.
- Han, A. K. (1987). Non-parametric analysis of a generalized regression model: the maximum rank correlation estimator. *Journal of Econometrics*, 35(2–3), 303–316.
- Härdle, W., Hall, P., & Ichimura, H. (1993). Optimal smoothing in single-index models. *The Annals of Statistics*, 21(1), 157–178.
- Hooke, R., & Jeeves, T. A. (1961). “Direct search” solution of numerical and statistical problems. *Journal of the ACM*, 8(2), 212–229.
- Hristache, M., Juditsky, A., & Spokoiny, V. (2001). Direct estimation of the index coefficient in a single-index model. *The Annals of Statistics*, 595–623.
- Huang, J. (1996). Efficient estimation for the proportional hazards model with interval censoring. *The Annals of Statistics*, 24(2), 540–568.
- Huang, J., & Wellner, J. (1993). Regression models with interval censoring. *Proceedings of the Kolmogorov Seminar, Euler Mathematics Institute*, St. Petersburg, Russia.
- Ichimura, H. (1993). Semiparametric least squares (SLS) and weighted SLS estimation of single-index models. *Journal of Econometrics*, 58(1–2), 71–120.
- Karabatsos, G. (2009). Modeling heteroscedasticity in the single-index model with the Dirichlet process. *Advances and Applications in Statistical Sciences*, 1(1), 83–104.
- Klein, R. W., & Spady, R. H. (1993). An efficient semiparametric estimator for binary response models. *Econometrica*, 61(2), 387–421.
- Kuchibhotla, A. K., & Patra, R. K. (2017). Efficient estimation in convex single index models. Retrieved from <https://arxiv.org/abs/1708.00145>
- Li, G., & Zhang, C.-H. (1998). Linear regression with interval censored data. *The Annals of Statistics*, 26(4), 1306–1327.
- Manski, C. F. (1975). Maximum score estimation of the stochastic utility model of choice. *Journal of Econometrics*, 3(3), 205–228.
- Murphy, S. A., van der Vaart, A. W., & Wellner, J. A. (1999). Current status regression. *Mathematical Methods of Statistics*, 8(3), 407–425.
- Newey, W. K., & Stoker, T. M. (1993). Efficiency of weighted average derivative estimators and index models. *Econometrica: Journal of the Econometric Society*, 1199–1223.

- Robertson, T., Wright, F., & Dykstra, R. (1988). *Order restricted statistical inference*, Wiley Series in Probability and Mathematical Statistics: Probability and Mathematical Statistics. Chichester: John Wiley & Sons.
- Sherman, R. P. (1993). The limiting distribution of the maximum rank correlation estimator. *Econometrica*, 61(1), 123–137.
- Torczon, V. (1997). On the convergence of pattern search algorithms. *SIAM Journal on Optimization*, 7(1), 1–25.
- Wang, H.-B. (2009). Bayesian estimation and variable selection for single index models. *Computational Statistics & Data Analysis*, 53(7), 2617–2627.
- Xia, Y., & Härdle, W. (2006). Semi-parametric estimation of partially linear single-index models. *Journal of Multivariate Analysis*, 97(5), 1162–1184.
- Yu, Y., & Ruppert, D. (2002). Penalized spline estimation for partially linear single-index models. *Journal of the American Statistical Association*, 97(460), 1042–1054.

How to cite this article: Groeneboom P, Hendrickx K. Estimation in monotone single-index models. *Statistica Neerlandica*. 2019;73:78–99. <https://doi.org/10.1111/stan.12138>

APPENDIX

A.1 | Equivalence MRCE and MRE for the binary choice model and the current status regression model

In this section, we show that the maximizers of H_n and R_n coincide for the current status linear regression model. A similar derivation holds for the binary choice model.

Consider the MRE for the current status model maximizing

$$R_n(\alpha) = \sum_{i \neq j} \Delta_i \{T_i - \alpha^T \mathbf{X}_i \beta > T_j - \alpha^T \mathbf{X}_j\}.$$

This rank estimator is equivalent to Han's MRCE, given by the maximizer of

$$\begin{aligned} H_n(\alpha) &= \sum_{i \neq j} \{\Delta_i > \Delta_j\} \{T_i - \alpha^T \mathbf{X}_i \beta > T_j - \alpha^T \mathbf{X}_j\} \\ &= \sum_{i \neq j} \Delta_i (1 - \Delta_j) \{T_i - \alpha^T \mathbf{X}_i \beta > T_j - \alpha^T \mathbf{X}_j\}. \end{aligned}$$

This can be seen as follows. Suppose that the observations are ordered in the $T_i - \alpha^T \mathbf{X}_i$, that is, $T_1 - \alpha^T \mathbf{X}_1 \leq T_2 - \alpha^T \mathbf{X}_2 \leq \dots \leq T_n - \alpha^T \mathbf{X}_n$. Then,

$$\begin{aligned} H_n(\alpha) &= \sum_{j < i} \Delta_i (1 - \Delta_j) = \sum_{j < i} \Delta_i - \sum_{j < i} \Delta_i \Delta_j = \sum_{j < i} \Delta_i - \frac{1}{2} \sum_{j \neq i} \Delta_i \Delta_j \\ &= R_n(\alpha) - \frac{1}{2} \sum_{j \neq i} \Delta_i \Delta_j. \end{aligned}$$

Since the second term in the expression above is independent of the ordering in $T_i - \alpha^T \mathbf{X}_i$, the maximizers of $R_n(\alpha)$ and $H_n(\alpha)$ coincide and both estimators are equivalent.

A.2 | Asymptotic distribution of the MRCE and the MRE

In this section, we show that the asymptotic normal distribution for the MRCE is given by

$$\sqrt{n}(\hat{\alpha}_n - \alpha_0) \rightarrow_d N_d(\mathbf{0}, \mathbf{V}^- \Sigma, \mathbf{V}^-),$$

where \mathbf{V} and Σ are defined in Table 1 and where \mathbf{V}^- denotes the Moore–Penrose inverse of matrix \mathbf{V} . A similar argument can be used to derive the asymptotic distribution of the MRE in terms of Moore–Penrose inverses. The asymptotic normality for the MRCE and the MRE is derived in Sherman (1993) and Cavanagh and Sherman (1998), where the authors restrict the parameter space to a compact subset of $\{\alpha \in \mathbb{R}^d : \alpha_d = 1\}$. Each α is represented as $(\theta, 1)$, and only $d - 1$ instead of d components are considered in the proofs of \sqrt{n} -consistency and asymptotic normality. Using the parametrization $\{\alpha \in \mathbb{R}^d : \|\alpha\| = 1\}$ instead and considering a transformation $\mathbb{S} : B \subset \mathbb{R}^{d-1} \mapsto \{\alpha \in \mathbb{R}^d : \|\alpha\| = 1\}$ as in Balabdaoui et al. (2017), it follows, by similar arguments as in Sherman, that for the MRCE, we have

$$\begin{aligned} 0 \leq H_n(\beta) - H_n(\beta_0) &= \frac{1}{2}(\beta - \beta_0)^T \mathbb{E} \{ \nabla_2 \tau(\mathbf{X}, Y), \beta \} (\beta - \beta_0) + \frac{1}{\sqrt{n}}(\beta - \beta_0)^T \mathbf{W}_n \\ &+ o_p(\|\beta - \beta_0\|^2) + o_p(1/n), \end{aligned}$$

where

$$\tau(\mathbf{x}, y, \beta) = \mathbb{E}(\{y > Y\} \{\mathbb{S}(\beta)^T \mathbf{x} > \mathbb{S}(\beta)^T \mathbf{X}\}) + \mathbb{E}(\{Y > y\} \{\mathbb{S}(\beta)^T \mathbf{X} > \mathbb{S}(\beta)^T \mathbf{x}\})$$

and

$$\mathbf{W}_n := \frac{1}{\sqrt{n}} \sum_{i=1}^n \nabla_1 \tau(\mathbf{X}_i, Y_i, \beta_0)$$

converges in distribution to a normal $N_{d-1}(\mathbf{0}, \mathbf{W})$ random vector. Here, ∇_1 represents the first partial derivative operator with respect to β . Let \mathbf{A} denote $\mathbb{E}(\nabla_2 \tau(\mathbf{X}, Y), \beta)$. If \mathbf{A} is negative definite, then it follows, by theorem 2 in Sherman, that

$$\sqrt{n}(\hat{\beta}_n - \beta_0) \rightarrow_d N_{d-1}(\mathbf{0}, \mathbf{A}^{-1} \mathbf{W} \mathbf{A}^{-1}).$$

Using an application of the delta method, we conclude that

$$\sqrt{n}(\hat{\alpha}_n - \alpha_0) \rightarrow_d N_d(\mathbf{0}, [\nabla_1 \mathbb{S}(\beta_0)] [\mathbf{A}^{-1} \mathbf{W} \mathbf{A}^{-1}] [\nabla_1 \mathbb{S}(\beta_0)]^T) = (\mathbf{0}, \mathbf{V}^- \Sigma \mathbf{V}^-),$$

where \mathbf{V} and Σ are given in Table 1 and where the last equality follows analogously to the proof of asymptotic normality given in Balabdaoui et al. (2017).