



J. R. Statist. Soc. B (2017)
79, Part 4, pp. 1165–1185

On estimation of optimal treatment regimes for maximizing t -year survival probability

Runchao Jiang, Wenbin Lu, Rui Song and Marie Davidian

North Carolina State University, Raleigh, USA

[Received October 2014. Revised June 2016]

Summary. A treatment regime is a deterministic function that dictates personalized treatment based on patients' individual prognostic information. There is increasing interest in finding optimal treatment regimes, which determine treatment at one or more treatment decision points to maximize expected long-term clinical outcomes, where larger outcomes are preferred. For chronic diseases such as cancer or human immunodeficiency virus infection, survival time is often the outcome of interest, and the goal is to select treatment to maximize survival probability. We propose two non-parametric estimators for the survival function of patients following a given treatment regime involving one or more decisions, i.e. the so-called value. On the basis of data from a clinical or observational study, we estimate an optimal regime by maximizing these estimators for the value over a prespecified class of regimes. Because the value function is very jagged, we introduce kernel smoothing within the estimator to improve performance. Asymptotic properties of the proposed estimators of value functions are established under suitable regularity conditions, and simulation studies evaluate the finite sample performance of the regime estimators. The methods are illustrated by application to data from an acquired immune deficiency syndrome clinical trial.

Keywords: Inverse probability weighted estimation; Kaplan–Meier estimator; Optimal treatment regime; Personalized medicine; Survival probability; Value function

1. Introduction

For many complex diseases, such as cancer, human immunodeficiency virus infection and mental disorders, there is generally not a uniformly best treatment for all patients. Rather, different patients may benefit from different treatments due to individual heterogeneity. For example, in AIDS Clinical Trials Group (ACTG) study 175 (Hammer *et al.*, 1996), the primary composite outcome of interest was time to having a larger than 50% decline in CD4 cell count, which is a measure of immunological status; progression to acquired immune deficiency syndrome or death. For the comparison of two treatments, zidovudine plus didanosine (coded as 1) and zidovudine plus zalcitabine (coded as 0), the data suggest that zidovudine plus zalcitabine leads to more favourable outcomes for younger patients than zidovudine plus didanosine. Fig. 1 shows treatment-specific Kaplan–Meier estimates of the survival function for the two age strata defined by the observed median age, 34 years, in ACTG study 175. It is clear that, among younger patients with age 34 years or less, those receiving zidovudine plus zalcitabine have almost uniformly larger survival probabilities than those receiving zidovudine plus didanosine, whereas the situation is reversed for older patients with age more than 34 years.

This type of situation suggests that individual patient characteristics should be used when

Address for correspondence: Wenbin Lu, Department of Statistics, North Carolina State University, 5212 SAS Hall, 2311 Stinson Drive, Raleigh, NC 27695, USA.
E-mail: lu@stat.ncsu.edu

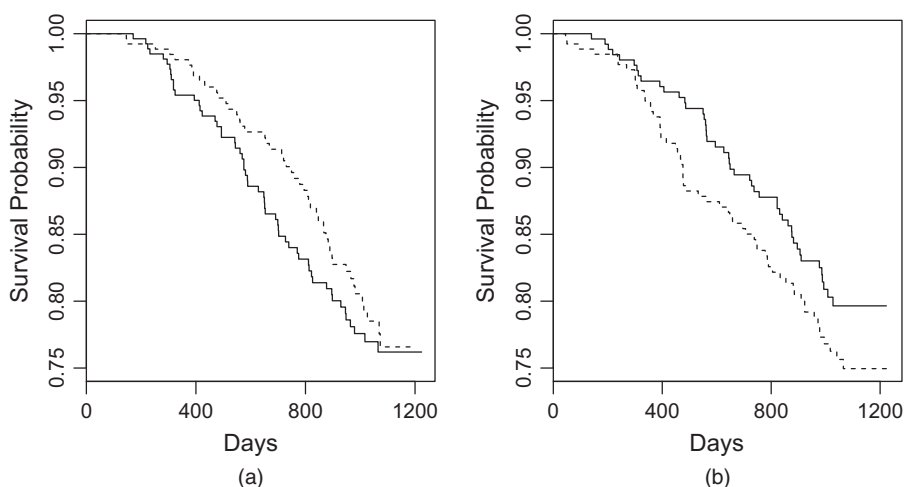


Fig. 1. Treatment-specific Kaplan–Meier curves by age (—, treatment 1; ----, treatment 0): (a) age 34 years or less; (b) age more than 34 years

selecting treatments to maximize an expected long-term outcome of interest for which larger outcomes are preferred, such as the t -year survival probability, and has heightened interest in derivation of optimal dynamic treatment regimes. Because in many chronic diseases treatment decisions may be made sequentially over time, a dynamic treatment regime is a set of one or more decision rules to determine which treatment to give from among the available options based on accruing individual patient information, including baseline characteristics, intermediate outcomes between decisions and previous treatments. An optimal regime is a regime that maximizes the expected outcome, or so-called value, if used by the entire patient population to select treatments.

There is a large literature on statistical methods to estimate an optimal treatment regime based on data from a clinical trial or observational study and non-survival outcomes. Q -learning (Watkins, 1989; Watkins and Dayan, 1992; Murphy, 2005; Zhao *et al.*, 2009) and A -learning (Murphy, 2003; Robins, 2004) are two popular backward induction methods for estimating optimal dynamic treatment regimes based on regression-type modelling. The former involves positing parametric models for, roughly, the regression of outcome on accruing information and treatment, whereas the latter is based on semiparametric models in which only the part of the outcome regression representing contrasts between treatments is modelled parametrically, along with the propensity scores, which are the probabilities of observed treatment assignment given patient information at each decision point. Q -learning can be sensitive to misspecification of the required models, whereas A -learning enjoys the so-called double-robustness property in that the corresponding estimating equations are asymptotically unbiased when either the propensity scores or main effects portion of the outcome models are correctly specified. An alternative class of approaches known as value or policy search methods is based on deriving and maximizing directly a consistent estimator for the value over a prespecified class of treatment regimes indexed by a finite dimensional parameter. Zhang *et al.* (2012b) proposed inverse propensity score weighted and augmented inverse propensity score weighted estimators for the value in the case of a single decision point. Because the value estimator is non-smooth, the optimization problem is challenging, and non-standard optimization techniques are required. Zhao *et al.* (2012) and Zhang *et al.* (2012a) recast this approach as a weighted classification problem; the former referred to this method as outcome-weighted learning. These approaches exploit approximations

integrated into classification software to address the non-smooth optimization problem, so that the class of regimes is dictated by a chosen classification method. Zhang *et al.* (2013) extended the value search methods of Zhang *et al.* (2012b) to more than one decision point, which share the computational challenges in the single-decision case. Matsouaka *et al.* (2014) employed a kernel smoothing technique to estimate non-parametrically the conditional mean for the difference of the potential outcomes in a subgroup of patients and derived its associated treatment regime.

Although the survival time is often the outcome of interest, to our knowledge there is relatively little development of methods for estimation of optimal treatment regimes where the goal is to maximize the survival probability. Some work has focused on maximizing the expected survival time. Goldberg and Kosorok (2012) developed a Q -learning method for censored survival data for estimating optimal dynamic treatment regimes and derived its associated finite sample risk bounds on the generalization error of the estimated regime, whereas Zhao *et al.* (2015) proposed a doubly robust estimator for expected survival time based on censored data and used outcome-weighted learning to estimate an optimal regime. Bai *et al.* (2014) developed a locally efficient doubly robust estimator for survival probability rather than mean survival time and estimated an optimal regime by extending the methods from a classification perspective of Zhang *et al.* (2012a). The last two methods involve transforming maximization of the value to a weighted classification problem, which allows classification software to be used to address the optimization challenge and thus dictates the class of regimes. All these methods are relevant to a single decision point only.

In this paper, we propose a value search method for estimating an optimal treatment regime within a prespecified class for which the goal is to maximize the survival probability that addresses the optimization challenges in a novel way and is relevant to more than one decision point. In particular, we develop a framework employing kernel smoothing techniques to smooth the estimator of the value before optimization, which we show greatly improves the finite sample performance over the corresponding estimator with no smoothing. This approach is different from the smoothing technique that was used by Matsouaka *et al.* (2014), and, to the best of our knowledge, this is the first time smoothing has been integrated into estimation of the value function and its associated optimal treatment regimes in this way. Development of optimal treatment regimes for multiple decision points with censored survival data is challenging, as timing of observations, censoring and events must be properly taken into account. In addition, we extend our smoothing approach to this setting.

In Sections 2 and 3, we introduce the statistical framework and estimators for a single decision point and multiple decisions respectively. Asymptotic properties of the estimators proposed are given in Section 4. Finite sample performance is studied via simulation in Section 5, and Section 6 presents application of the methods to data from the ACTG study 175. Proofs are relegated to Appendix A.

The data that are analysed in the paper and the programs that were used to analyse them can be obtained from

<http://wileyonlinelibrary.com/journal/rss-datasets>

2. Estimation of optimal treatment regime for a single decision time point

2.1. Notation and assumptions

Consider a study with two treatment options $\mathcal{A} = \{0, 1\}$ given at baseline. For the i th patient, $i = 1, \dots, n$, let \mathbf{X}_i denote the p -dimensional vector of baseline covariates taking values $\mathbf{x} \in \mathcal{X}$ and A_i denote the actual treatment received by the patient. Let T_i be the associated continuous survival time of interest, with conditional survival function $S_T(t|a, \mathbf{x}) \equiv P(T_i > t|A_i = a, \mathbf{X}_i = \mathbf{x})$

and corresponding conditional cumulative hazard function $\Lambda_T(t|a, \mathbf{x})$, where $a = 0, 1$. Let C_i denote the right censoring time for patient i . The observed data are $\{(\mathbf{X}_i, A_i, \tilde{T}_i, \delta_i), i = 1, \dots, n\}$, independent and identically distributed across i , where $\tilde{T}_i = \min\{T_i, C_i\}$ and $\delta_i = I(T_i \leq C_i)$. We thus observe the counting process $N_i(t) = I(\tilde{T}_i \leq t, \delta_i = 1)$ and the at-risk process $Y_i(t) = I(\tilde{T}_i \geq t)$.

A treatment regime is a deterministic function that maps $\mathbf{x} \in \mathcal{X}$ to \mathcal{A} . For simplicity, we assume that the regimes of interest are from $\mathcal{G} = \{g_\eta : g_\eta(\mathbf{x}) = I(\boldsymbol{\eta}^T \tilde{\mathbf{x}} \geq 0), \|\boldsymbol{\eta}\| = 1\}$, where $\tilde{\mathbf{x}} = (1, \mathbf{x}^T)^T$. However, the method proposed also applies to any other \mathcal{G} indexed by finite dimensional parameters. Denote the potential survival time of a patient if he or she were given treatment a , which may be contrary to fact, as $T^*(a)$. Accordingly, define the potential counting process $N^*(a; t)$ and the at-risk process $Y^*(a; t)$ under treatment a , where $N^*(a; t) = I[\min\{T^*(a), C\} \leq t, T^*(a) \leq C]$ and $Y^*(a; t) = I[\min\{T^*(a), C\} \geq t]$. If a patient follows a given regime g_η , we can write the corresponding potential survival time as $T^*(g_\eta) = T^*(1)g_\eta + T^*(0)(1 - g_\eta)$, whose survival function is given by $S^*(t; \boldsymbol{\eta}) = E[P\{T^*\{g_\eta(\mathbf{X})\} > t | \mathbf{X}\}]$, as well as the potential counting process $N^*(g_\eta; t) = N^*(1; t)g_\eta + N^*(0; t)(1 - g_\eta)$ and potential at-risk process $Y^*(g_\eta; t) = Y^*(1; t)g_\eta + Y^*(0; t)(1 - g_\eta)$. We wish to find an optimal treatment regime in \mathcal{G} that maximizes t -year survival probability, i.e. $g_\eta^{\text{opt}}(\mathbf{x}) \equiv g(\mathbf{x}; \boldsymbol{\eta}^{\text{opt}})$, where $\boldsymbol{\eta}^{\text{opt}} = \arg \max_{\|\boldsymbol{\eta}\|=1} S^*(t; \boldsymbol{\eta})$. Here, t is a predetermined time point.

To find an optimal treatment regime, we first derive consistent estimators of $S^*(u; \boldsymbol{\eta})$ for any u . We make the uninformative censoring assumption $\{T^*(1), T^*(0)\} \perp\!\!\!\perp C | A, \mathbf{X}$, where ' $\perp\!\!\!\perp$ ' means 'independent of'. Let $S_C(t|a, \mathbf{x})$ denote the survival function of the censoring time given $A = a$ and $\mathbf{X} = \mathbf{x}$. If we could observe the g_η -specific potential counting process $N_i^*(g_\eta; s)$ and the at-risk process $Y_i^*(g_\eta; s)$, an intuitive estimator for $S^*(u; \boldsymbol{\eta})$ is the inverse probability of censoring weighted Kaplan–Meier estimator

$$\hat{S}^*(u; \boldsymbol{\eta}) = \prod_{s \leq u} \left(1 - \frac{\sum_{i=1}^n [dN_i^*\{g_\eta(\mathbf{X}_i); s\} / S_C\{s | g_\eta(\mathbf{X}_i), \mathbf{X}_i\}]}{\sum_{i=1}^n [Y_i^*\{g_\eta(\mathbf{X}_i); s\} / S_C\{s | g_\eta(\mathbf{X}_i), \mathbf{X}_i\}]} \right). \quad (1)$$

However, because $N_i^*(g_\eta; s)$ and $Y_i^*(g_\eta; s)$ are generally not observable, $\hat{S}^*(u; \boldsymbol{\eta})$ is not computable on the basis of the observed data. To obtain proper estimators that are computable from the observed data, we make the following two assumptions, which are widely used in the causal inference literature (Rubin, 1974):

- (a) the stable unit treatment value assumption, i.e. $T = T^*(1)A + T^*(0)(1 - A)$, and
- (b) the no unmeasured confounders assumptions, i.e. $\{T^*(1), T^*(0)\} \perp\!\!\!\perp A | \mathbf{X}$.

2.2. Estimation procedure

Following Zhang *et al.* (2012b), we cast estimation of $S^*(u; \boldsymbol{\eta})$ in a missing data framework. By the stable unit treatment value assumption, for those patients whose actually received treatment matches the treatment dictated by g_η , $N_i^*(g_\eta; s) = N_i(s)$ and $Y_i^*(g_\eta; s) = Y_i(s)$, which are observed. For other patients, they are missing. This motivates us to modify the estimator that is given in equation (1) by incorporating inverse propensity score weighting. Formally, the weight for the i th patient is given by

$$\begin{aligned} w_{\eta i} &= \frac{I\{A_i = I(\boldsymbol{\eta}^T \tilde{\mathbf{X}} \geq 0)\}}{\pi(\mathbf{X}_i)A_i + \{1 - \pi(\mathbf{X}_i)\}(1 - A_i)} \\ &= \frac{A_i I(\boldsymbol{\eta}^T \tilde{\mathbf{X}} \geq 0) + (1 - A_i)\{1 - I(\boldsymbol{\eta}^T \tilde{\mathbf{X}} \geq 0)\}}{\pi(\mathbf{X}_i)A_i + \{1 - \pi(\mathbf{X}_i)\}(1 - A_i)}, \end{aligned} \quad (2)$$

where $\pi(\mathbf{X}_i) = P(A_i = 1 | \mathbf{X}_i)$ is the propensity score. In practice, $\pi(\mathbf{X}_i)$ is known by design, as in a randomized clinical trial, or must be modelled and estimated from the data as in observational studies. In the latter case, a parametric model, say a logistic regression, is usually used for estimating $\pi(\mathbf{X}_i)$: specifically,

$$\text{logit}\{\pi(\mathbf{X}_i; \boldsymbol{\theta})\} = \boldsymbol{\theta}^T \tilde{\mathbf{X}}_i, \quad (3)$$

where $\text{logit}(z) = \log\{z/(1-z)\}$. Let $\hat{\boldsymbol{\theta}}$ denote the maximum likelihood estimator of $\boldsymbol{\theta}$, and define $\hat{\pi}(\mathbf{X}_i) = \exp(\hat{\boldsymbol{\theta}}^T \tilde{\mathbf{X}}_i) / \{1 + \exp(\hat{\boldsymbol{\theta}}^T \tilde{\mathbf{X}}_i)\}$. If the logistic regression model is correctly specified, $\hat{\boldsymbol{\theta}}$ is a consistent estimator of $\boldsymbol{\theta}$.

To derive the estimator for $S^*(u; \boldsymbol{\eta})$, we need to estimate the censoring time survival function $S_C(s | A_i, \mathbf{X}_i)$ also. In many clinical studies with satisfactory follow-up, it is reasonable to assume that censoring times are independent of treatment assignment and covariates, i.e. independent censoring. Here, the Kaplan–Meier estimator for censoring times consistently estimates $S_C(s | A_i, \mathbf{X}_i)$. For some applications, the independent censoring assumption may be restrictive but can be relaxed to a certain extent. For example, if censoring times are assumed to depend only on treatment assignment, the stratified Kaplan–Meier estimator can be used to estimate the treatment-specific censoring time survival function. For more general dependence, we can build a semiparametric model, say a proportional hazards model for censoring times, and obtain the model-based estimator of $S_C(s | A_i, \mathbf{X}_i)$. For simplicity, from now on we make the independent censoring assumption and let $\hat{S}_C(\cdot)$ denote the Kaplan–Meier estimator for censoring times.

Let $\hat{w}_{\eta i}$ denote the estimator of $w_{\eta i}$ that is obtained by replacing $\pi(\mathbf{X}_i)$ with $\hat{\pi}(\mathbf{X}_i)$ in $w_{\eta i}$. We propose the inverse propensity score weighted Kaplan–Meier estimator (IPSWKME) for $S^*(u; \boldsymbol{\eta})$ given by

$$\hat{S}_I(u; \boldsymbol{\eta}) = \prod_{s \leq u} \left\{ 1 - \frac{\sum_{i=1}^n \hat{w}_{\eta i} dN_i(s)}{\sum_{i=1}^n \hat{w}_{\eta i} Y_i(s)} \right\}. \quad (4)$$

Note that the IPSWKME does not depend on the Kaplan–Meier estimator $\hat{S}_C(\cdot)$ for censoring times, as it cancels from the numerator and denominator under the independent censoring assumption. In Section 4, we show that $\hat{S}_I(u; \boldsymbol{\eta})$ is a consistent estimator of $S^*(u; \boldsymbol{\eta})$ under certain conditions. Based on $\hat{S}_I(u; \boldsymbol{\eta})$, the estimated optimal treatment regime to maximize the t -year survival probability is given by $g(\mathbf{x}; \hat{\boldsymbol{\eta}}_I^{\text{opt}})$, where $\hat{\boldsymbol{\eta}}_I^{\text{opt}} = \arg \max_{\|\boldsymbol{\eta}\|=1} \hat{S}_I(t; \boldsymbol{\eta})$.

The IPSWKME (4) relies on a correct specification of the propensity score model. If it is misspecified, the IPSWKME is inconsistent. To improve the robustness of the IPSWKME, we propose the augmented IPSWKME (AIPSWKME) by incorporating assumed model information. For example, we may posit a proportional hazards model (Cox, 1972) for the conditional cumulative hazard function of T by

$$\Lambda_T(t | A, \mathbf{X}) = \Lambda_0(t) \exp\{\boldsymbol{\beta}^T (\mathbf{X}^T, A, A\mathbf{X}^T)^T\}, \quad (5)$$

where $\Lambda_0(t)$ is the baseline cumulative hazard function, and $\boldsymbol{\beta}$ is a $(2p+1)$ -dimensional parameter. The term $w_{\eta i} dN_i^* \{g_{\eta}(\mathbf{X}_i); s\}$ is augmented by

$$\begin{aligned} & w_{\eta i} dN_i^* \{g_{\eta}(\mathbf{X}_i); s\} + (1 - w_{\eta i}) E[dN_i^* \{g_{\eta}(\mathbf{X}_i); s\} | \mathbf{X}_i] \\ & = w_{\eta i} dN_i^* \{g_{\eta}(\mathbf{X}_i); s\} + (1 - w_{\eta i}) S_T\{s | g_{\eta}(\mathbf{X}_i), \mathbf{X}_i\} S_C(s) d\Lambda_T\{s | g_{\eta}(\mathbf{X}_i), \mathbf{X}_i\}, \end{aligned}$$

where $S_T(s | A_i, \mathbf{X}_i)$ and $S_C(s)$ are the conditional survival functions of T and C respectively. Similarly, the term $w_{\eta i} Y_i^* \{g_{\eta}(\mathbf{X}_i); s\}$ is augmented by $w_{\eta i} Y_i^* \{g_{\eta}(\mathbf{X}_i); s\} + (1 - w_{\eta i}) S_T\{s | g_{\eta}(\mathbf{X}_i), \mathbf{X}_i\} S_C(s)$. It can be shown that the two augmented terms have the so-called double-robustness

property, i.e. they are unbiased for $E[dN_i^*\{g_\eta(\mathbf{X}_i); s\}|\mathbf{X}_i]$ and $E[Y_i^*\{g_\eta(\mathbf{X}_i); s\}|\mathbf{X}_i]$ respectively, when either the propensity score model or the posited proportional hazards model is correctly specified. Therefore, we propose the AIPSWKME for $S^*(u; \eta)$ as

$$\hat{S}_A(u; \eta) = \prod_{s \leq u} \left(1 - \frac{\sum_{i=1}^n [\hat{w}_{\eta_i} dN_i(s) + (1 - \hat{w}_{\eta_i}) \hat{S}_T\{s|g_\eta(\mathbf{X}_i), \mathbf{X}_i\} \hat{S}_C(s) d\hat{\Lambda}_T\{s|g_\eta(\mathbf{X}_i), \mathbf{X}_i\}]}{\sum_{i=1}^n [\hat{w}_{\eta_i} Y_i(s) + (1 - \hat{w}_{\eta_i}) \hat{S}_T\{s|g_\eta(\mathbf{X}_i), \mathbf{X}_i\} \hat{S}_C(s)]} \right), \quad (6)$$

where $\hat{S}_T(s|A_i, \mathbf{X}_i)$ is the estimated survival function of T based on the fitted proportional hazards model and $\hat{S}_C(s)$ is the Kaplan–Meier estimator for censoring times. Based on $\hat{S}_A(u; \eta)$, the estimated optimal treatment regime to maximize the t -year survival probability is given by $g(\mathbf{x}; \hat{\eta}_A^{\text{opt}})$, where $\hat{\eta}_A^{\text{opt}} = \arg \max_{\|\eta\|=1} \hat{S}_A(t; \eta)$. The asymptotic properties of $\hat{S}_A(u; \eta)$ and $\hat{S}_A(t; \hat{\eta}_A^{\text{opt}})$ are studied in Section 4.

2.3. Computational aspects

$\hat{S}_1(t; \eta)$ and $\hat{S}_A(t; \eta)$ are not smooth functions of η . As an illustration, we plot $\hat{S}_1(t; \eta)$ and $\hat{S}_A(t; \eta)$ as functions of η_1 in Fig. 2 for a simple example with one covariate and the intercept term in η set as 1. The estimates are very jagged, and direct maximization of them with respect to η is challenging and may lead to local maximizers. From our simulation studies in Section 5, estimated survival probabilities following the optimal treatment regimes obtained may show substantial biases. As studied in Matsouaka *et al.* (2014), cross-validation may be used to correct the finite sample biases of the unsmoothed estimators, but it may increase the computational burden.

To reduce the biases of the estimators, we propose to smooth the estimators $\hat{S}_1(t; \eta)$ and $\hat{S}_A(t; \eta)$ by using kernel smoothers. Specifically, we replace $g_\eta(\mathbf{x}_i) = I(\eta^T \tilde{\mathbf{x}}_i \geq 0)$ in $\hat{S}_1(t; \eta)$ and $\hat{S}_A(t; \eta)$ with $\tilde{g}_\eta(\mathbf{x}_i) = \Phi(\eta^T \tilde{\mathbf{x}}_i/h)$ to obtain the smoothed IPSWKME (SIPSWKME) $\tilde{S}_1(t; \eta)$ and smoothed AIPSWKME (SAIPSWKME) $\tilde{S}_A(t; \eta)$, where $\Phi(s)$ is the cumulative distribution function for the standard normal distribution, and h is a bandwidth parameter that goes to 0 as $n \rightarrow \infty$. For bandwidth selection, we set $h = c_0 n^{-1/3} \text{sd}(\eta^T \tilde{\mathbf{X}})$, where c_0 is a constant and $\text{sd}(\mathbf{v})$ is the sample standard deviation of \mathbf{v} . In our numerical studies, we found that $c_0 = 4^{1/3}$ generally gives good results for all scenarios. We plot in Fig. 2 the smoothed estimates with the chosen bandwidth parameter for the same example. The smoothed estimates approximate the original estimates well and have unique maximizers around the true value $\eta_1 = 0.5$. Because the treatment regime $I(\eta^T \tilde{\mathbf{X}} \geq 0)$ remains the same when η is multiplied by k for any $k > 0$, choosing the bandwidth h to be a function of η , in particular h being proportional to $\text{sd}(\eta^T \tilde{\mathbf{X}})$, ensures the scale-free property of the regime, as the constant k cancels in $\Phi(\eta^T \tilde{\mathbf{X}}/h)$. As shown in Fig. 2, although the resulting smoothed value function is not convex in η , it generally has a unique mode, and the maximizer of the smoothed value function is much easier to obtain compared with the unsmoothed counterpart. In all our numerical studies, the non-convexity of the smoothed value function does not cause any difficulty in the maximization procedure. Such a bandwidth parameter has been widely used in the non-parametric smoothing literature and ensures that the original and smoothed estimators have the same asymptotic distribution (e.g. Heller (2007)). Let $\tilde{\eta}_1^{\text{opt}}$ and $\tilde{\eta}_A^{\text{opt}}$ denote the maximizers of $\tilde{S}_1(t; \eta)$ and $\tilde{S}_A(t; \eta)$ respectively. Then the associated estimated optimal treatment regimes are $g(\mathbf{x}; \tilde{\eta}_1^{\text{opt}})$ and $g(\mathbf{x}; \tilde{\eta}_A^{\text{opt}})$. In our implementation, we first conduct the optimization without the norm 1 constraint. Instead, we search the maximizer in the domain $-1 \leq \eta_j \leq 1$ for all j s and then we rescale η to have norm 1. This does not change the estimated value function, \hat{S}_1 and \hat{S}_A , and their smoothed counterparts.

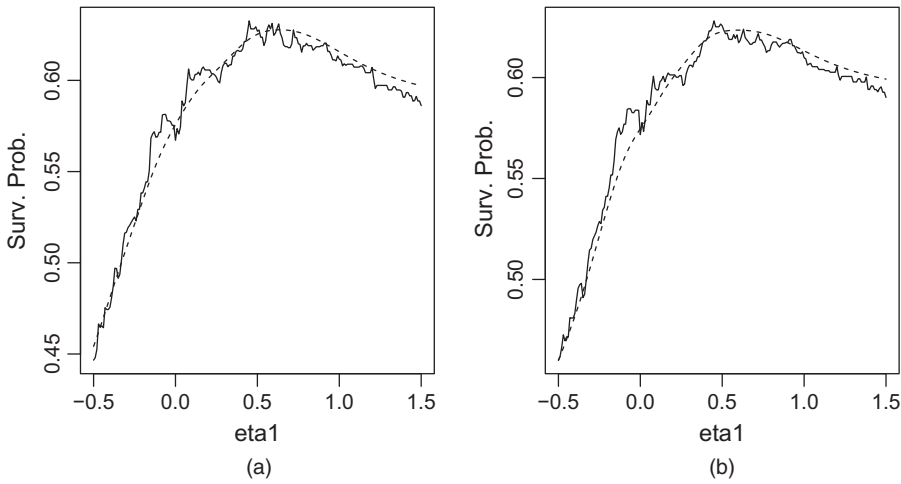


Fig. 2. Plots for the original (—) and smoothed (-----) estimates: (a) IPSWKME estimates; (b) AIP-SWKME estimates

3. Estimation of optimal treatment regime for multiple decision time points

We now extend the foregoing methods to estimation of optimal dynamic treatment regimes incorporating multiple decision points. For simplicity, we illustrate for the case of two decision points. Specifically, treatments can be given at baseline and at a fixed interim time point s , $0 < s < t$. For the i th patient, let \mathbf{X}_{0i} denote his or her p_0 -dimensional vector of baseline covariates and $A_{0i} \in \mathcal{A}_0 = \{0, 1\}$ denote the initial treatment received at baseline. If the patient survives beyond s and is not censored before s , let \mathbf{X}_{1i} denote his or her p_1 -dimensional vector of intermediate covariates collected by s after assigning treatment A_{0i} and $A_{1i} \in \mathcal{A}_1 = \{0, 1\}$ denote the follow-up treatment given at s . Thus, the observed data are $\{\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i} I(\tilde{T}_i > s), A_{1i} I(\tilde{T}_i > s), \tilde{T}_i, \delta_i, i = 1, \dots, n\}$ and independent and identically distributed across i .

As for a single decision point, we consider a class of linear dynamic treatment regimes for simplicity, i.e. $\mathcal{G} = \{\mathbf{g}_\eta = (g_0, g_1)\}$, where

$$g_0(\mathbf{x}_0; \boldsymbol{\eta}_0) = I\{\boldsymbol{\eta}_0^\top (1, \mathbf{x}_0^\top) \geq 0\},$$

$$g_1(\mathbf{x}_0, \mathbf{x}_1; \boldsymbol{\eta}_1) = I[\boldsymbol{\eta}_1^\top \{1, \mathbf{x}_0^\top, g_0(\mathbf{x}_0; \boldsymbol{\eta}_0), \mathbf{x}_1^\top\} \geq 0],$$

$\boldsymbol{\eta}_0$ is a $(p_0 + 1)$ -dimensional parameter with $\|\boldsymbol{\eta}_0\| = 1$ and $\boldsymbol{\eta}_1$ is a $(p_0 + p_1 + 2)$ -dimensional parameter with $\|\boldsymbol{\eta}_1\| = 1$. Here, a patient following a treatment regime \mathbf{g}_η is given treatment $g_0(\mathbf{X}_0; \boldsymbol{\eta}_0)$ at baseline, and, if he or she survives beyond s and is not censored before s , is given treatment $g_1(\mathbf{X}_0, \mathbf{X}_1; \boldsymbol{\eta}_1)$ at s . For patients whose initial treatments coincide with those assigned by $g_0(\mathbf{X}_0; \boldsymbol{\eta}_0)$ and who die before s , their treatment assignments are consistent with the regime \mathbf{g}_η . However, for patients whose initial treatments coincide with those assigned by $g_0(\mathbf{X}_0; \boldsymbol{\eta}_0)$ but who are censored before s , it is not known whether their treatment assignments at the second decision follow the regime \mathbf{g}_η . Let $T^*(\mathbf{g}_\eta)$ denote the potential survival time for a patient if he or she were given treatment according to $\mathbf{g}_\eta(\mathbf{X}_0, \mathbf{X}_1)$. We are interested in finding the optimal dynamic treatment regime $\mathbf{g}_\eta^{\text{opt}} = \{g_0(\mathbf{X}_0; \boldsymbol{\eta}_0^{\text{opt}}), g_1(\mathbf{X}_0, \mathbf{X}_1; \boldsymbol{\eta}_1^{\text{opt}})\}$ in \mathcal{G} that maximizes the t -year survival probability $S^{*(2)}(t; \boldsymbol{\eta}) = E(P[T^*\{\mathbf{g}_\eta(\mathbf{X}_0, \mathbf{X}_1)\} > t | \mathbf{X}_0, \mathbf{X}_1])$. As is standard in the causal inference literature for studying dynamic treatment regimes (e.g. Murphy (2003)), we assume

- (a) the stable unit treatment value assumption, i.e. a patient's observed outcome agrees with

the corresponding potential outcome if his or her actually received treatments are consistent with the assigned treatments, and

- (b) the sequential randomization assumption, i.e. the treatment assignment at the current stage depends on only the past received treatments and observed covariates, but not the potential outcomes.

Under these two assumptions, the above-defined t -year survival probability can be estimated from the observed data.

We propose a similar IPSWKME for the survival function $S^{*(2)}(u; \boldsymbol{\eta})$ given any treatment regime $\mathbf{g}_{\boldsymbol{\eta}}$. However, the derivation of proper weights becomes more difficult, as some patients may be censored before s and whether their received treatments follow regime $\mathbf{g}_{\boldsymbol{\eta}}$ is unknown. To take this into account, we define the following new weight for patient i , $i = 1, \dots, n$:

$$\hat{w}_{\boldsymbol{\eta}i}^{(2)} = \frac{I(\tilde{T}_i \leq s) \delta_i}{\hat{S}_C(\tilde{T}_i)} \frac{I\{A_{0i} = g_0(\mathbf{X}_{0i}; \boldsymbol{\eta}_0)\}}{\hat{\pi}_{A_0}(\mathbf{X}_{0i})} + \frac{I(\tilde{T}_i > s)}{\hat{S}_C(s)} \frac{I[A_{0i} = g_0(\mathbf{X}_{0i}; \boldsymbol{\eta}_0), A_{1i} = g_1\{\mathbf{X}_{0i}, g_0(\mathbf{X}_{0i}; \boldsymbol{\eta}_0), X_{1i}; \boldsymbol{\eta}_1\}]}{\hat{\pi}_{A_0}(\mathbf{X}_{0i}) \hat{\pi}_{A_1}(\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i})},$$

where $\hat{\pi}_{A_0}(\mathbf{X}_{0i}) = \hat{\pi}_0(\mathbf{X}_{0i}) A_{0i} + \{1 - \hat{\pi}_0(\mathbf{X}_{0i})\}(1 - A_{0i})$, $\hat{\pi}_{A_1}(\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i}) = \hat{\pi}_1(\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i}) A_{1i} + \{1 - \hat{\pi}_1(\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i})\}(1 - A_{1i})$, and $\hat{\pi}_0(\mathbf{X}_{0i})$ and $\hat{\pi}_1(\mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i})$ are the estimates of the propensity scores $P(A_{0i} = 1 | \mathbf{X}_{0i})$ and $P(A_{1i} = 1 | \mathbf{X}_{0i}, A_{0i}, \mathbf{X}_{1i}, \tilde{T}_i > s)$ respectively. In randomized studies, $\hat{\pi}_0$ and $\hat{\pi}_1$ are known by design, whereas, in observational studies, they must be estimated, e.g. by using logistic regression. The IPSWKME for $S^*(u; \boldsymbol{\eta})$ is given by

$$\hat{S}_I^{(2)}(u; \boldsymbol{\eta}) = \prod_{v \leq u} \left\{ 1 - \frac{\sum_{i=1}^n \hat{w}_{\boldsymbol{\eta}i}^{(2)} dN_i(v)}{\sum_{i=1}^n \hat{w}_{\boldsymbol{\eta}i}^{(2)} Y_i(v)} \right\}. \quad (7)$$

Let $\hat{\boldsymbol{\eta}}_I^{\text{opt},(2)} = (\hat{\boldsymbol{\eta}}_{1,0}^{\text{opt},(2)}, \hat{\boldsymbol{\eta}}_{1,1}^{\text{opt},(2)}) = \arg \max_{\|\boldsymbol{\eta}_0\|=1, \|\boldsymbol{\eta}_1\|=1} \hat{S}_I^{(2)}(t; \boldsymbol{\eta})$. Then the estimated optimal dynamic treatment regime is given by $\hat{\mathbf{g}}_{\boldsymbol{\eta}}^{\text{opt},(2)} = \{g_0(\mathbf{X}_0; \hat{\boldsymbol{\eta}}_{1,0}^{\text{opt},(2)}), g_1(\mathbf{X}_0, \mathbf{X}_1; \hat{\boldsymbol{\eta}}_{1,1}^{\text{opt},(2)})\}$.

To improve the finite sample performance of the IPSWKME, we again introduce kernel smoothing and replace the indicator functions $g_0(\mathbf{X}_{0i}; \boldsymbol{\eta}_0)$ and $g_1(\mathbf{X}_{0i}, \mathbf{X}_{1i}; \boldsymbol{\eta}_1)$ in $\hat{S}_I^{(2)}(u; \boldsymbol{\eta})$ by $\Phi\{\boldsymbol{\eta}_0^T(1, \mathbf{X}_{0i}^T)/h_0\}$ and $\Phi[\boldsymbol{\eta}_1^T\{1, \mathbf{X}_0^T, g_0(\mathbf{X}_0; \boldsymbol{\eta}_0), \mathbf{X}_1^T\}/h_1]$, where the bandwidth parameters h_0 and h_1 are chosen as before. Let $\tilde{S}_I^{(2)}(u; \boldsymbol{\eta})$ denote the resulting smoothed IPSWKME and $\tilde{\boldsymbol{\eta}}_I^{\text{opt},(2)}$ denote the maximizer of $\tilde{S}_I^{(2)}(t; \boldsymbol{\eta})$. To improve the robustness of the IPSWKME, we can similarly derive the AIPSWKME based on a posited model for survival time; however, its formulation is very complicated and is not pursued here. Conceptually, the IPSWKME proposed can be generalized to accommodate more than two decision points. However, when there are more treatment decision points, the IPSWKME may become less reliable because fewer patients will have assigned treatments that are consistent with a given dynamic treatment regime.

4. Asymptotic properties

In this section, we present the asymptotic properties of the proposed estimators in theorems 1–3. Theorems 1 and 2 are for the cases with a single decision point whereas theorem 3 is for the case with two decision points.

Theorem 1. Under conditions 1–6 in Appendix A, if the propensity score model (3) is correctly specified, for any regime $\mathbf{g}_{\boldsymbol{\eta}}$, we have, as $n \rightarrow \infty$,

- (a) $\hat{S}_I(u; \boldsymbol{\eta}) \rightarrow^P S^*(u; \boldsymbol{\eta})$ for any $0 < u \leq t$,
- (b) $\sqrt{n}\{\hat{S}_I(u; \boldsymbol{\eta}) - S^*(u; \boldsymbol{\eta})\}$ converges weakly to a mean 0 Gaussian process,
- (c) $\sqrt{n}\{\hat{S}_I(t; \hat{\boldsymbol{\eta}}_I^{\text{opt}}) - S^*(t; \boldsymbol{\eta}^{\text{opt}})\} \rightarrow^d N\{0, \Sigma_I(t; \boldsymbol{\eta}^{\text{opt}})\}$, where the expression for $\Sigma_I(t; \boldsymbol{\eta}^{\text{opt}})$ is given in Appendix A, and
- (d) $\sqrt{n}\{\hat{S}_I(t; \hat{\boldsymbol{\eta}}_I^{\text{opt}}) - \tilde{S}_I(t; \tilde{\boldsymbol{\eta}}_I^{\text{opt}})\} = o_p(1)$.

Theorem 2. Under conditions 1–6 in Appendix A, if either the propensity score model (3) or the proportional hazard model (5) is correctly specified, we have, as $n \rightarrow \infty$,

- (a) $\hat{S}_A(u; \boldsymbol{\eta}) \rightarrow^P S^*(u; \boldsymbol{\eta})$ for any $0 < u \leq t$,
- (b) $\sqrt{n}\{\hat{S}_A(u; \boldsymbol{\eta}) - S^*(u; \boldsymbol{\eta})\}$ converges weakly to a mean 0 Gaussian process,
- (c) $\sqrt{n}\{\hat{S}_A(t; \hat{\boldsymbol{\eta}}_A^{\text{opt}}) - S^*(t; \boldsymbol{\eta}^{\text{opt}})\} \rightarrow^d N\{0, \Sigma_A(t; \boldsymbol{\eta}^{\text{opt}})\}$, where the expression for $\Sigma_A(t; \boldsymbol{\eta}^{\text{opt}})$ is given in Appendix A, and
- (d) $\sqrt{n}\{\hat{S}_A(t; \hat{\boldsymbol{\eta}}_A^{\text{opt}}) - \tilde{S}_A(t; \tilde{\boldsymbol{\eta}}_A^{\text{opt}})\} = o_p(1)$.

Theorem 3. Under certain regularity conditions, if the two propensity score models $\pi_0(\cdot)$ and $\pi_1(\cdot)$ are correctly specified, for any regime $g_{\boldsymbol{\eta}}$, we have, as $n \rightarrow \infty$,

- (a) $\hat{S}_I^{(2)}(u; \boldsymbol{\eta}) \rightarrow^P S^{*(2)}(u; \boldsymbol{\eta})$ for any $0 < u \leq t$,
- (b) $\sqrt{n}\{\hat{S}_I^{(2)}(u; \boldsymbol{\eta}) - S^{*(2)}(u; \boldsymbol{\eta})\}$ converges weakly to a mean 0 Gaussian process,
- (c) $\sqrt{n}\{\hat{S}_I^{(2)}(t; \hat{\boldsymbol{\eta}}_I^{\text{opt}, (2)}) - S^{*(2)}(t; \boldsymbol{\eta}^{\text{opt}, (2)})\} \rightarrow^d N\{0, \Sigma_I^{(2)}(t; \boldsymbol{\eta}^{\text{opt}, (2)})\}$, where $\boldsymbol{\eta}^{\text{opt}, (2)} = (\boldsymbol{\eta}_0^{\text{opt}}, \boldsymbol{\eta}_1^{\text{opt}})$, and
- (d) $\sqrt{n}\{\hat{S}_I^{(2)}(t; \hat{\boldsymbol{\eta}}_I^{\text{opt}, (2)}) - \tilde{S}_I^{(2)}(t; \tilde{\boldsymbol{\eta}}_I^{\text{opt}, (2)})\} = o_p(1)$.

Here the asymptotic variances $\Sigma_I(t; \boldsymbol{\eta}^{\text{opt}})$, $\Sigma_A(t; \boldsymbol{\eta}^{\text{opt}})$ and $\Sigma_I^{(2)}(t; \boldsymbol{\eta}^{\text{opt}, (2)})$ can be consistently estimated from the observed data by using the usual plug-in method. The proofs of theorems 1–3 are given in Appendix A.

5. Simulation studies

We examine the finite sample performance of the proposed estimators by simulation. We first consider scenarios with a single treatment decision time point at baseline. For each patient, baseline covariates X_1 and X_2 are independently and uniformly distributed on $(-2, 2)$. Given X_1 and X_2 , the binary treatment indicator A is generated from the logistic model $\text{logit}\{\pi(X_1, X_2)\} = X_1 - 0.5X_2$. The survival time T is generated from a linear transformation model (Cheng *et al.*, 1995), $h(T) = -0.5X_1 + A(X_1 - X_2) + \varepsilon$, where $h(s) = \log\{\exp(s) - 1\} - 2$ is an increasing function, and the error term ε follows some known distribution, either the extreme value distribution or the logistic distribution, which corresponds to a proportional hazards and proportional odds model respectively. The covariate-independent censoring time C is uniformly distributed on $(0, C_0)$, where C_0 is chosen to achieve a censoring rate of 15% and 40%. The optimal treatment regime maximizing t -year survival probability is $g_{\boldsymbol{\eta}}^{\text{opt}}(X_1, X_2) = I(X_1 - X_2 \geq 0)$ for any t . We search the optimal treatment regime in the class of regimes given by $\mathcal{G} = \{g_{\boldsymbol{\eta}} : g_{\boldsymbol{\eta}}(X_1, X_2) = I(\eta_0 + \eta_1 X_1 + \eta_2 X_2 \geq 0), \boldsymbol{\eta} = (\eta_0, \eta_1, \eta_2)^T\}$, which contains the true optimal treatment regime with $\boldsymbol{\eta}^{\text{opt}} = (0, 0.707, -0.707)$.

To implement the estimators, it is necessary to posit a model for the propensity scores. We consider both a correctly specified model, $\text{logit}\{\pi_A(X_1, X_2)\} = \theta_0 + \theta_1 X_1 + \theta_2 X_2$, and a misspecified model, $\text{logit}\{\pi_A(X_1, X_2)\} = \theta_0$. For the augmented estimators, we must posit a model for the survival time T . Here, we always use the proportional hazard model $\lambda(t|X_1, X_2) = \lambda_0(t) \exp\{\beta_{11} X_1 + \beta_{12} X_2 + A(\beta_{20} + \beta_{21} X_1 + \beta_{22} X_2)\}$. Note that, when ε follows the extreme value distribution, the posited survival model is correctly specified. In contrast, when ε follows

the logistic distribution, this model is misspecified. We compare the performance of the IPSWKME \hat{S}_I and AIPSWKME \hat{S}_A , as well as their smoothed versions the SIPSWKME \tilde{S}_I and SAIPSWKME \tilde{S}_A , under different combinations of the assumed propensity score model, error term distribution, censoring rate, sample size ($n = 250$ or $n = 500$) and time point of interest ($t = 1$ or $t = 2$). For each scenario, we ran 1000 replications and used a genetic algorithm to do the optimization, which is implemented by the R function `genoud` within the package `rgenoud` (Mebane and Sekhon, 2011).

We report results for the scenarios with $n = 250$ and $t = 2$, which are given in Tables 1 and 2 for the extreme value error and logistic error distributions respectively. Results for other scenarios are similar. In Tables 1 and 2 we report the mean of estimated η^{opt} , the mean of the estimated t -year survival probability following the estimated optimal treatment regime, namely the estimated optimal t -year survival probability (denoted by $\hat{S}(\hat{\eta}^{\text{opt}})$), the mean of the estimated standard error of $\hat{S}(\hat{\eta}^{\text{opt}})$ by using the plug-in method based on the asymptotic variances that were established in theorems 1 and 2 (denoted by SE), the empirical coverage probability of the 95% confidence interval for the t -year survival probability following the true optimal treatment regime $S(\eta^{\text{opt}})$ (denoted by CP), the mean of the simulated true t -year survival probability following the estimated optimal treatment regime (denoted by $S(\hat{\eta}^{\text{opt}})$) and the mean of the misclassification rate by comparing the true and estimated optimal treatment regimes (denoted by MR). The numbers given in parentheses are the standard deviations of the corresponding estimates. Here, $S(\eta^{\text{opt}})$ and $S(\hat{\eta}^{\text{opt}})$ are computed by using simulated survival times following the given treatment regime based on a large random sample of 5×10^6 patients. We have $S(\eta^{\text{opt}}) = 0.605$ for the extreme value error distribution and $S(\eta^{\text{opt}}) = 0.672$ for the logistic distribution. The misclassification rate for one simulation is calculated as the proportion of patients for which the true and estimated optimal treatment regimes do not select the same treatment.

From the results, when the propensity score model is correctly specified, all estimators of η^{opt} have relatively small biases; in particular, the mean of $\hat{\eta}_0^{\text{opt}}$ is close to 0 whereas the mean ratio of $\hat{\eta}_1^{\text{opt}}$ to $\hat{\eta}_2^{\text{opt}}$ is very close to -1 . The means of the simulated true t -year survival probability following the estimated optimal treatment regimes, i.e. $S(\hat{\eta}^{\text{opt}})$, are all close to the true values. In addition, the estimates of η^{opt} based on the AIPSWKME and SAIPSWKME of the t -year survival probability generally have smaller standard deviation than those based on the IPSWKME and SIPSWKME. The unsmoothed IPSWKME and AIPSWKME of the optimal t -year survival probability have relatively large biases mainly due to the very jagged estimates of the t -year survival probability, as illustrated in Fig. 2 and, as a consequence, the associated coverage probability of the 95% confidence interval is much lower than the nominal level. The SIPSWKME and SAIPSWKME of the optimal t -year survival probability greatly reduce the biases and thus give the proper coverage probability. In addition, the unsmoothed and smoothed estimators of the optimal t -year survival probability have nearly the same standard deviation. When the propensity score model is misspecified, the IPSWKME and SIPSWKME generally have relatively large biases as expected, whereas the AIPSWKME and SAIPSWKME greatly reduce the biases and give much smaller MR. In particular, when the posited survival model is correctly specified under the extreme value error distribution, the SAIPSWKME yields the proper coverage probability. In contrast, when the posited survival model is misspecified under the logistic error distribution, although the SAIPSWKME is not consistent in general, it still gives small biases with reasonable coverage probability. The performance of the estimators improves as the censoring rate decreases and the sample size increases.

We also compare the method proposed with the methods of Zhao *et al.* (2013, 2015). For the comparison with the method of Zhao *et al.* (2013), we consider randomized studies with known propensity scores, i.e. $\pi_A \equiv 0.5$, sample size $n = 250$, decision time point of interest $t_0 = 2$,

Table 1. Simulation results for the extreme value error distribution with $n = 250$ and $t = 2$

	PS^\dagger	$\hat{\eta}_0$	$\hat{\eta}_1$	$\hat{\eta}_2$	$\hat{S}(\hat{\eta}^{\text{opt}})^\ddagger$	SE	CP	$S(\hat{\eta}^{\text{opt}})^\ddagger$	MR
<i>Censor rate 15%</i>									
\hat{S}_I	T	0.010 (0.298)	0.633 (0.192)	-0.665 (0.178)	0.645 (0.037)	0.040	0.839	0.590 (0.016)	0.118 (0.063)
\tilde{S}_I	T	-0.005 (0.263)	0.652 (0.179)	-0.667 (0.171)	0.612 (0.036)	0.040	0.968	0.593 (0.014)	0.107 (0.057)
\hat{S}_A	T	-0.002 (0.287)	0.639 (0.171)	-0.676 (0.155)	0.639 (0.037)	0.040	0.866	0.592 (0.014)	0.109 (0.058)
\tilde{S}_A	T	0.002 (0.256)	0.654 (0.169)	-0.675 (0.152)	0.609 (0.036)	0.040	0.969	0.594 (0.013)	0.102 (0.055)
\hat{S}_I	F	-0.031 (0.423)	0.408 (0.327)	-0.697 (0.249)	0.666 (0.036)	0.039	0.659	0.565 (0.040)	0.193 (0.100)
\tilde{S}_I	F	-0.051 (0.403)	0.426 (0.285)	-0.714 (0.252)	0.643 (0.035)	0.039	0.844	0.569 (0.034)	0.184 (0.090)
\hat{S}_A	F	-0.014 (0.278)	0.660 (0.151)	-0.662 (0.161)	0.635 (0.038)	0.041	0.886	0.593 (0.012)	0.107 (0.055)
\tilde{S}_A	F	-0.002 (0.246)	0.675 (0.141)	-0.665 (0.148)	0.607 (0.038)	0.041	0.968	0.596 (0.010)	0.096 (0.050)
<i>Censor rate 40%</i>									
\hat{S}_I	T	0.008 (0.311)	0.616 (0.214)	-0.661 (0.202)	0.650 (0.041)	0.044	0.850	0.588 (0.019)	0.127 (0.068)
\tilde{S}_I	T	-0.002 (0.285)	0.637 (0.202)	-0.660 (0.191)	0.613 (0.040)	0.045	0.958	0.590 (0.017)	0.118 (0.064)
\hat{S}_A	T	0.006 (0.310)	0.623 (0.203)	-0.663 (0.189)	0.645 (0.041)	0.045	0.879	0.589 (0.019)	0.123 (0.068)
\tilde{S}_A	T	0.001 (0.282)	0.643 (0.192)	-0.661 (0.183)	0.612 (0.040)	0.044	0.965	0.591 (0.017)	0.115 (0.062)
\hat{S}_I	F	0.002 (0.448)	0.388 (0.349)	-0.676 (0.267)	0.671 (0.039)	0.043	0.677	0.560 (0.045)	0.206 (0.109)
\tilde{S}_I	F	-0.024 (0.432)	0.403 (0.311)	-0.694 (0.271)	0.645 (0.039)	0.043	0.867	0.564 (0.038)	0.200 (0.095)
\hat{S}_A	F	-0.005 (0.299)	0.655 (0.169)	-0.650 (0.176)	0.641 (0.043)	0.046	0.896	0.591 (0.014)	0.115 (0.060)
\tilde{S}_A	F	-0.005 (0.270)	0.664 (0.162)	-0.656 (0.173)	0.609 (0.041)	0.046	0.964	0.593 (0.012)	0.109 (0.054)

† PS, propensity score model; T means the correctly specified propensity score model whereas F means the misspecified propensity score model.
 ‡ Recall that $S(\eta^{\text{opt}}) = 0.605$.

Table 2. Simulation results for the logistic error distribution with $n = 250$ and $t = 2$

PS^\dagger	$\hat{\eta}_0$	$\hat{\eta}_1$	$\hat{\eta}_2$	$\hat{S}(\hat{\eta}^{opt})_+^\ddagger$	SE	CP	$S(\hat{\eta}^{opt})_+^\ddagger$	MR
Censor rate 15%								
\hat{S}_1	0.010 (0.370)	0.566 (0.272)	-0.641 (0.241)	0.716 (0.034)	0.038	0.791	0.652 (0.022)	0.155 (0.089)
\tilde{S}_1	-0.004 (0.341)	0.593 (0.262)	-0.640 (0.235)	0.685 (0.034)	0.039	0.955	0.655 (0.020)	0.145 (0.082)
\hat{S}_A	0.007 (0.363)	0.578 (0.260)	-0.639 (0.240)	0.713 (0.034)	0.039	0.818	0.653 (0.020)	0.151 (0.084)
\tilde{S}_A	-0.006 (0.341)	0.595 (0.251)	-0.642 (0.233)	0.684 (0.034)	0.039	0.962	0.655 (0.020)	0.143 (0.081)
\hat{S}_1	0.041 (0.461)	0.340 (0.389)	-0.662 (0.284)	0.729 (0.033)	0.037	0.649	0.632 (0.040)	0.224 (0.120)
\tilde{S}_1	-0.001 (0.461)	0.375 (0.350)	-0.667 (0.283)	0.707 (0.033)	0.037	0.846	0.636 (0.035)	0.216 (0.107)
\hat{S}_A	-0.025 (0.337)	0.630 (0.198)	-0.637 (0.210)	0.723 (0.036)	0.040	0.753	0.658 (0.013)	0.133 (0.068)
\tilde{S}_A	-0.029 (0.320)	0.633 (0.204)	-0.642 (0.204)	0.695 (0.036)	0.040	0.926	0.659 (0.012)	0.130 (0.064)
Censor rate 40%								
\hat{S}_1	0.013 (0.395)	0.545 (0.293)	-0.625 (0.266)	0.721 (0.036)	0.041	0.785	0.649 (0.027)	0.168 (0.097)
\tilde{S}_1	-0.008 (0.362)	0.581 (0.274)	-0.626 (0.255)	0.687 (0.036)	0.041	0.948	0.652 (0.022)	0.155 (0.087)
\hat{S}_A	0.004 (0.381)	0.558 (0.277)	-0.635 (0.255)	0.718 (0.036)	0.042	0.807	0.651 (0.023)	0.160 (0.089)
\tilde{S}_A	-0.016 (0.361)	0.578 (0.270)	-0.634 (0.246)	0.686 (0.036)	0.042	0.955	0.653 (0.022)	0.153 (0.086)
\hat{S}_1	0.061 (0.471)	0.325 (0.413)	-0.640 (0.299)	0.733 (0.035)	0.039	0.661	0.628 (0.042)	0.235 (0.124)
\tilde{S}_1	0.021 (0.482)	0.355 (0.370)	-0.639 (0.312)	0.709 (0.035)	0.040	0.842	0.631 (0.038)	0.229 (0.114)
\hat{S}_A	-0.012 (0.350)	0.625 (0.206)	-0.631 (0.217)	0.722 (0.038)	0.042	0.785	0.657 (0.014)	0.138 (0.070)
\tilde{S}_A	-0.022 (0.331)	0.628 (0.214)	-0.634 (0.221)	0.692 (0.038)	0.043	0.939	0.658 (0.013)	0.136 (0.067)

† PS, propensity score model; here T means the correctly specified propensity score model whereas F means the misspecified propensity score model.
 ‡ Recall that $S(\eta^{opt}) = 0.672$.

and a censoring rate of 15%. When implementing the method of Zhao *et al.* (2013), we set the threshold $\xi = 0, 0.1, \dots, 0.6$ and find the associated treatment regime for each ξ -value.

Table 3 summarizes the simulation results for the extreme value and logistic error distributions based on 1000 replications. The performance of the method of Zhao *et al.* (2013) depends on the choice of the threshold value ξ . For the extreme value error distribution, the best choice is $\xi = 0.4$, whereas, for the logistic error distribution, the best choice is $\xi = 0.3$. In practice, the best threshold value to use is unknown and must be estimated from data, which may not be straightforward. Moreover, even with the best choice of ξ -value, the performance of the method by Zhao *et al.* (2013) is still worse than that of our proposed smoothed estimators SIPSWKME and SAIPSWKME, under all the settings considered.

For the comparison with the method of Zhao *et al.* (2015), we consider the same simulation settings as in Tables 1 and 2 with sample size $n = 250$, decision time point of interest $t_0 = 2$ and a censoring rate of 15%. For both methods, we consider the augmented estimation. Table 4 summarizes the simulation results based on 1000 replications. The methods proposed and the method of Zhao *et al.* (2015) lead to comparable survival probabilities under the estimated treatment rules, whereas the methods proposed yield smaller misclassification rates under all the settings considered. In summary, the methods proposed demonstrate very competitive performance compared with existing approaches.

Next, we consider scenarios with two treatment decision time points: one at the baseline and the other at $s = 1$. The initial treatment assignment A_0 and the follow-up treatment assignment A_1 , if applicable, are generated independently from a Bernoulli distribution with success probability 0.5. A single baseline covariate X_0 is generated from a uniform distribution on $(0, 4)$. To generate the survival time T , we first generate a time T_1 given A_0 and X_0 from an exponential distribution with the rate function $\lambda_1(A_0, X_0)$. The censoring time C is generated from a uniform distribution on $(0, C_0)$. If a patient is neither dead nor censored at time $s = 1$ (i.e. $\min(T_1, C) > 1$), we generate a single intermediate covariate X_1 for this patient as $X_1 = 0.5X_0 - 0.4(A_0 - 0.5) + e$, where e is uniformly distributed on $(0, 2)$. Then we generate another time T_2 given A_0, A_1, X_0, X_1 from an exponential distribution with rate function $\lambda_2(A_0, A_1, X_0, X_1)$. The survival time T of interest is defined as $T = T_1$ if $T_1 \leq 1$ and $T = 1 + T_2$ otherwise. The observed survival time is $\tilde{T} = \min(T, C)$ with the censoring indicator $\delta = I(T \leq C)$. Here, C_0 is chosen to achieve censoring rates of 15% and 40%. We consider three scenarios for the rate functions λ_1 and λ_2 :

- (a) $\lambda_1(A_0, X_0) = 0.5 \exp\{1.75(A_0 - 0.5)(X_0 - 2)\}$ and $\lambda_2(A_0, A_1, X_0, X_1) = 0.3 \exp\{2.5(A_1 - 0.4)(X_1 - 2) - A_0(X_1 - 2)\}$;
- (b) $\lambda_1(A_0, X_0) = 0.1 \exp\{2(A_0 - 0.5)(X_0 - 2)\}$ and $\lambda_2(A_0, A_1, X_0, X_1) = 0.2 \exp\{3(A_1 - 0.4) \times (X_1 - 2) - 3(A_0 - 0.5)(X_0 - 2)\}$;
- (c) $\lambda_1(A_0, X_0) = 0.2 \exp\{1.5(A_0 - 0.3)(X_0 - 3)\}$ and $\lambda_2(A_0, A_1, X_0, X_1) = 0.3 \exp\{2(A_1 - 0.5) \times (X_1 - 2) + 0.5(A_0 - 0.7)(X_0 - 1)\}$.

For these three scenarios, the true optimal rule for maximizing t -year survival probability ($t > 1$) at time $s = 1$ is given by $g_1^{\text{opt}}(x_1) = I(2 - x_1 > 0)$. However, the true optimal rule $g_0^{\text{opt}}(x_0)$ at time $s = 0$ is a complicated non-linear function of x_0 , which can be derived by using backward induction as in Q -learning. In our implementation, for computational simplicity, we search for the optimal dynamic treatment regime in a class involving linear decision rules: specifically, $\mathcal{G}_\eta = \{g_0(x_0) = I(\eta_1 + \eta_2 x_0 > 0), g_1(x_1) = I(\eta_3 + \eta_4 x_1 > 0), \|(\eta_1, \eta_2)\| = 1, \|(\eta_3, \eta_4)\| = 1\}$. Then, the true optimal rule $g_1^{\text{opt}}(x_1)$ at time $s = 1$ corresponds to $(\eta_3^{\text{opt}}, \eta_4^{\text{opt}}) = (0.894, -0.447)$ for all three scenarios.

For scenarios (a) and (c), we take $t = 3$, whereas for (b) we take $t = 6$. We use simulation to find the true optimal rule at $s = 0$ in \mathcal{G}_η to maximize the t -year survival probability. Specifically, we first

Table 3. Results for comparison with the method of Zhao *et al.* (2013)†

Error	Method	Survival probability	MR
Extreme value	Zhao <i>et al.</i> (2013) with $\xi = 0$	0.445 (0.030)	0.467 (0.048)
	Zhao <i>et al.</i> (2013) with $\xi = 0.1$	0.499 (0.046)	0.373 (0.089)
	Zhao <i>et al.</i> (2013) with $\xi = 0.2$	0.555 (0.035)	0.245 (0.099)
	Zhao <i>et al.</i> (2013) with $\xi = 0.3$	0.585 (0.027)	0.143 (0.091)
	Zhao <i>et al.</i> (2013) with $\xi = 0.4$	0.590 (0.028)	0.112 (0.093)
	Zhao <i>et al.</i> (2013) with $\xi = 0.5$	0.543 (0.066)	0.241 (0.162)
	Zhao <i>et al.</i> (2013) with $\xi = 0.6$	0.542 (0.045)	0.275 (0.109)
	SIPSWKME	0.594 (0.011)	0.107 (0.052)
Logistic	SAIPSWKME	0.595 (0.009)	0.099 (0.047)
	Zhao <i>et al.</i> (2013) with $\xi = 0$	0.552 (0.028)	0.456 (0.061)
	Zhao <i>et al.</i> (2013) with $\xi = 0.1$	0.606 (0.040)	0.323 (0.113)
	Zhao <i>et al.</i> (2013) with $\xi = 0.2$	0.643 (0.029)	0.200 (0.111)
	Zhao <i>et al.</i> (2013) with $\xi = 0.3$	0.650 (0.030)	0.164 (0.117)
	Zhao <i>et al.</i> (2013) with $\xi = 0.4$	0.630 (0.037)	0.246 (0.132)
	Zhao <i>et al.</i> (2013) with $\xi = 0.5$	0.590 (0.039)	0.373 (0.110)
	Zhao <i>et al.</i> (2013) with $\xi = 0.6$	0.590 (0.030)	0.382 (0.079)
	SIPSWKME	0.659 (0.012)	0.130 (0.063)
	SAIPSWKME	0.660 (0.011)	0.126 (0.061)

†Survival probability, the simulated survival probability at $t_0 = 2$; MR, the misclassification rate. The true optimal survival probabilities are 0.605 and 0.672 for the extreme value and logistic error respectively. Values in parentheses are the standard deviations over 1000 simulations.

Table 4. Results for comparison with the method of Zhao *et al.* (2015)†

Error	Method	PS	Survival probability	MR
Extreme value	Zhao <i>et al.</i> (2015)	T	0.587 (0.022)	0.136 (0.065)
	AIPSWKM	T	0.592 (0.014)	0.109 (0.058)
	SAIPSWKM	T	0.594 (0.013)	0.102 (0.055)
	Zhao <i>et al.</i> (2015)	F	0.590 (0.008)	0.134 (0.044)
	AIPSWKM	F	0.593 (0.012)	0.107 (0.055)
	SAIPSWKM	F	0.596 (0.010)	0.096 (0.050)
Logistic	Zhao <i>et al.</i> (2015)	T	0.652 (0.027)	0.159 (0.090)
	AIPSWKM	T	0.653 (0.020)	0.151 (0.084)
	SAIPSWKM	T	0.655 (0.020)	0.143 (0.081)
	Zhao <i>et al.</i> (2015)	F	0.659 (0.007)	0.141 (0.047)
	AIPSWKM	F	0.658 (0.013)	0.133 (0.068)
	SAIPSWKM	F	0.659 (0.012)	0.130 (0.064)

†PS, propensity score model: here T means the correctly specified propensity score model whereas F means the misspecified propensity score model. Survival probability, the simulated survival probability at $t_0 = 2$; MR, the misclassification rate. The true optimal survival probabilities are 0.605 and 0.672 for the extreme value and logistic error respectively. Values in parentheses are the standard deviations over 1000 simulations.

generate X_0 and, for a given (η_1, η_2) , we set A_0 by the regime $g_0(X_0)$. Then, we generate X_1 given A_0 and X_0 the same way as in our design and set A_1 by the optimal rule g_1^{opt} . Finally, we generate T_1 and T_2 , and define T the same way as before. Using generated T s for a large random sample of 5×10^6 patients, we compute the associated empirical t -year survival probability. We find $(\eta_1^{\text{opt}}, \eta_2^{\text{opt}})$ to maximize the empirical t -year survival probability, which gives the true optimal rule g_0^{opt} in \mathcal{G}_η . Here, we use the grid search method to find $(\eta_1^{\text{opt}}, \eta_2^{\text{opt}})$. Since

Table 5. Simulation results for estimating optimal dynamic treatment regimes[†]

C%	S	η_1^{opt}	η_2^{opt}	η_3^{opt}	η_4^{opt}	$\hat{S}(\hat{\eta}^{\text{opt}})$	SE	CP	$S(\hat{\eta}^{\text{opt}})$	MR
<i>Scenario 1: $\eta^{\text{opt}} = (0.890, -0.456, 0.894, -0.447)$; $S(3; \eta^{\text{opt}}) = 0.567$</i>										
15	F	0.881 (0.035)	-0.468 (0.062)	0.893 (0.017)	-0.449 (0.033)	0.591 (0.028)	0.030	0.887	0.559 (0.008)	0.107 (0.054)
	T	0.884 (0.029)	-0.463 (0.052)	0.894 (0.013)	-0.448 (0.026)	0.570 (0.028)	0.030	0.955	0.561 (0.006)	0.088 (0.048)
40	F	0.878 (0.042)	-0.471 (0.072)	0.890 (0.022)	-0.453 (0.042)	0.600 (0.036)	0.037	0.841	0.555 (0.011)	0.125 (0.061)
	T	0.884 (0.034)	-0.463 (0.061)	0.892 (0.018)	-0.450 (0.035)	0.574 (0.035)	0.038	0.955	0.558 (0.009)	0.108 (0.056)
<i>Scenario 2: $\eta^{\text{opt}} = (-0.891, 0.454, 0.894, -0.447)$; $S(6; \eta^{\text{opt}}) = 0.624$</i>										
15	F	-0.888 (0.025)	0.456 (0.045)	0.891 (0.018)	-0.452 (0.035)	0.645 (0.025)	0.027	0.890	0.615 (0.008)	0.099 (0.052)
	T	-0.889 (0.018)	0.456 (0.034)	0.893 (0.014)	-0.450 (0.028)	0.624 (0.024)	0.027	0.967	0.618 (0.005)	0.079 (0.042)
40	F	-0.886 (0.030)	0.459 (0.053)	0.890 (0.020)	-0.453 (0.038)	0.650 (0.027)	0.029	0.855	0.613 (0.010)	0.110 (0.055)
	T	-0.888 (0.022)	0.457 (0.040)	0.892 (0.016)	-0.451 (0.032)	0.626 (0.027)	0.030	0.972	0.617 (0.007)	0.091 (0.048)
<i>Scenario 3: $\eta^{\text{opt}} = (0.908, -0.419, 0.894, -0.447)$; $S(3; \eta^{\text{opt}}) = 0.702$</i>										
15	F	0.897 (0.037)	-0.434 (0.069)	0.892 (0.020)	-0.449 (0.039)	0.728 (0.026)	0.027	0.829	0.692 (0.009)	0.134 (0.067)
	T	0.900 (0.031)	-0.430 (0.060)	0.893 (0.016)	-0.448 (0.031)	0.707 (0.026)	0.027	0.952	0.695 (0.007)	0.116 (0.061)
40	F	0.895 (0.041)	-0.437 (0.075)	0.891 (0.023)	-0.451 (0.043)	0.732 (0.028)	0.029	0.809	0.691 (0.010)	0.142 (0.073)
	T	0.899 (0.035)	-0.431 (0.066)	0.893 (0.019)	-0.449 (0.036)	0.709 (0.028)	0.030	0.951	0.693 (0.008)	0.126 (0.066)

[†] C% denotes the censoring rate; \hat{S} indicates whether the smoothing technique is applied (T) or not (F).

$\|(\eta_1^{\text{opt}}, \eta_2^{\text{opt}})\| = 1$, we need to do grid search only for η_1 . We have $(\eta_1^{\text{opt}}, \eta_2^{\text{opt}}) = (0.890, -0.456)$ and $S(3; \eta^{\text{opt}}) = 0.567$ for scenario 1, $(\eta_1^{\text{opt}}, \eta_2^{\text{opt}}) = (-0.891, 0.454)$ and $S(6; \eta^{\text{opt}}) = 0.624$ for scenario 2 and $(\eta_1^{\text{opt}}, \eta_2^{\text{opt}}) = (0.908, -0.419)$ and $S(3; \eta^{\text{opt}}) = 0.702$ for scenario 3. Here $\eta^{\text{opt}} = (\eta_1^{\text{opt}}, \eta_2^{\text{opt}}, \eta_3^{\text{opt}}, \eta_4^{\text{opt}})^T$ and $S(t; \eta^{\text{opt}})$ is the t -year survival probability following the optimal dynamic treatment regime that is defined by η^{opt} .

We compare the unsmoothed and smoothed estimators. For both estimators, the propensity score models π_0 and π_1 are assumed known as for randomized clinical trials. Simulation results for 1000 replications are summarized in Table 5. From the results, we observe that

- both unsmoothed and smoothed estimation methods give nearly unbiased estimators of η^{opt} , and the t -year survival probability following the estimated optimal treatment regime (which is denoted by $S(\hat{\eta}^{\text{opt}})$ in Table 5) is very close to the t -year survival probability following the true optimal treatment regime η^{opt} ,
- the mean of the estimated standard error SE of $\hat{S}(\hat{\eta}^{\text{opt}})$ based on the established theory is close to the standard deviation of the estimates given in parentheses,
- the unsmoothed estimator for the t -year survival probability following the estimated optimal treatment regime (which is denoted by $\hat{S}(\hat{\eta}^{\text{opt}})$) has relatively large bias and the associated coverage probability CP is below the nominal level and
- the smoothed estimator for the t -year survival probability following the estimated optimal treatment regime has largely reduced bias and thus leads to the proper coverage probability.

6. Application to AIDS Clinical Trials Group study 175

We illustrate the proposed methods for a single decision with the data from the ACTG study 175 (Hammer *et al.*, 1996). Subjects were randomized to four treatment groups with equal probability: zidovudine monotherapy, ZDV, ZDV plus didanosine, ddI, ZDV plus zalcitabine, zal, and ddI monotherapy. A primary composite end point of interest is the time to having a larger than 50% decline in the CD4 cell count, or progressing to acquired immune deficiency syndrome or death, whichever comes first. From treatment-specific Kaplan–Meier curves, it can be clearly seen that treatments ZDV + ddI, ZDV + zal and ddI only are uniformly better than treatment ZDV only in terms of survival. In addition, treatments ZDV + ddI and ZDV + zal are overall the two best treatments giving the highest survival probabilities especially after day 400. For simplicity, we consider only two treatment options in our analysis: $A = 1$ for ZDV + ddI and $A = 0$ for ZDV + zal, which involves 1046 subjects. For each subject, there are 12 baseline clinical covariates; preliminary analysis results showed that Karnofsky score Karnof, baseline CD4 cell count CD40 and age Age are three important risk predictors and may have interaction effects with treatments. We include these three covariates in constructing treatment regimes. Our goal is to find the optimal treatment regime from the class of linear regimes defined by $\mathcal{G} = \{g_\eta = I(\eta_0 + \eta_1 x_1 + \eta_2 x_2 + \eta_3 x_3 \geq 0) : \eta = (\eta_0, \eta_1, \eta_2, \eta_3)^T, \|\eta\| = 1\}$ to maximize t -year survival probability, x_1 is Karnof, x_2 is CD40 and x_3 is Age. Because the data come from a randomized study, we use a constant model for the propensity score and estimate this constant from data. For augmented estimation, we posit the proportional hazard model as given in expression (5). We estimate optimal treatment regimes at day $t = 400, 600, 800, 1000$.

The estimated optimal treatment regimes and the associated t -year survival probabilities are presented in Table 6. We present only the results for the SIPSWKME and SAIPSWKME, as they have better numerical performance than their non-smoothed counterparts on the basis of our simulation studies. The numbers in the columns of Intercept, Karnof, CD40 and Age are the parameter estimates $\hat{\eta}^{\text{opt}}$ defining the optimal treatment regimes, and $\tilde{S}(t; \hat{\eta}^{\text{opt}})$ is the estimated



Table 6. Estimation results for the ACTG study 175 data†

t	Method	Intercept	Karnof	CD40	Age	$\tilde{S}(t; \tilde{\eta}^{\text{opt}})$	CI_1	CI_0
400	I	−0.303	−0.340	0.024	0.890	0.965 (0.008)	(−0.002, 0.023)	(−0.003, 0.044)
	A	−0.729	−0.240	0.018	0.640	0.965 (0.008)	(−0.002, 0.022)	(−0.003, 0.043)
600	I	0.975	−0.082	0.001	0.206	0.923 (0.012)	(0.000, 0.045)	(−0.006, 0.052)
	A	0.909	−0.137	0.000	0.392	0.922 (0.012)	(0.000, 0.043)	(−0.009, 0.052)
800	I	0.871	−0.133	−0.010	0.473	0.887 (0.014)	(0.008, 0.058)	(−0.002, 0.069)
	A	0.874	−0.131	−0.009	0.469	0.886 (0.014)	(0.006, 0.057)	(−0.003, 0.068)
1000	I	−0.210	−0.185	−0.035	0.959	0.824 (0.017)	(0.004, 0.060)	(−0.006, 0.081)
	A	0.001	−0.187	−0.037	0.982	0.823 (0.017)	(0.002, 0.059)	(−0.007, 0.080)

†I denotes the SIPSWKME and A denotes the SAIPSWKME; the numbers in parentheses are the estimated standard errors; CI_1 and CI_0 denote the 95% confidence intervals for the difference of the value functions obtained under the estimated optimal treatment regime and the simple treatment regime assigning all to treatment 1 and 0 respectively.

t -year survival probability following the estimated optimal treatment regime. From Table 6, the estimated optimal treatment regime for an earlier time may be different from that for a later time. For example, comparing the optimal treatment regimes obtained for $t = 600$ and $t = 800$, the SIPSWKME assigns a set of 353 patients to treatment 0 and another set of 583 patients to treatment 1 for both time points. However, it assigns a set of 52 patients to treatment 0 for $t = 600$ but to treatment 1 for $t = 800$. In contrast, it assigns another set of 58 patients to treatment 1 for $t = 600$ but to treatment 0 for $t = 800$. For the SAIPSWKME, the findings are similar. The SIPSWKME and SAIPSWKME yield very different parameter estimates $\tilde{\eta}^{\text{opt}}$. However, the corresponding optimal treatment regimes are similar. Using the results for day 600 as an example, among the 1046 subjects, there are only 57 subjects whose assigned treatments are different by the estimated optimal treatment regimes based on the SIPSWKME and SAIPSWKME. In addition, the estimated t -year survival probabilities following the estimated optimal treatment regimes are nearly the same on the basis of the SIPSWKME and SAIPSWKME.

Next, we compare the estimated optimal regimes with the simple regimes that assign all subjects to the same treatment. Specifically, we construct 95% Wald-type confidence intervals for the difference between the estimated t -year survival probabilities under the estimated optimal treatment regimes and the simple regimes based on the derived asymptotic normal distribution. The results are also given in Table 6. The confidence intervals either stay above 0 or 0 is very close to the left end point of the intervals when it is contained. This implies that the increase in value that is realized by following the estimated optimal treatment regimes comparing with the simple regimes is significant or at least marginally significant. The Kaplan–Meier curves for patients following the estimated optimal treatment regimes (not shown here) are all uniformly better than those for each single treatment.

We also estimated the optimal treatment regimes by using the proposed methods based on all 12 covariates when smoothing is and is not employed. We do not report on this here for brevity; however, we note that the results for the smoothed estimators when using three versus 12 covariates are comparable, demonstrating the adaptivity of the smoothed estimators to incorporating relatively many covariates. The unsmoothed estimators can lead to slightly different optimal treatment rules but with similar estimated survival probabilities. In addition, the estimated survival probabilities show relatively larger differences between the cases with three and 12 covariates, which is likely to be due to the instability in maximizing the unsmoothed value functions.



7. Discussion

We have proposed Kaplan–Meier-type estimators for the survival function of patients following a given (dynamic) treatment regime and introduced kernel smoothing to improve their performance. An optimal (dynamic) treatment regime within a class of prespecified treatment regimes may then be estimated by maximizing the estimator of the associated t -year survival probability. We consider the case when there are two treatment options at each decision time point. However, the methods proposed can be generalized to incorporate multiple treatment options at each decision by defining a treatment regime by using multiple indices instead of a single indicator function. In addition, current methods find the optimal (dynamic) treatment regime to maximize the t -year survival probability, which can also be generalized to maximize other clinical outcomes of interest. Specifically, using the IPSWKME, $\hat{S}_1(\cdot; \boldsymbol{\eta})$, as an illustration, we can find the optimal treatment regime to maximize $f\{\hat{S}_1(\cdot; \boldsymbol{\eta})\}$, where f is a specified function of interest; for example, $f\{\hat{S}_1(\cdot; \boldsymbol{\eta})\} = \int_0^L \hat{S}_1(u; \boldsymbol{\eta}) du$ corresponds to a restricted mean survival time under a given treatment regime. Likewise $f\{\hat{S}_1(\cdot; \boldsymbol{\eta})\} = \sup\{u : \hat{S}_1(u; \boldsymbol{\eta}) \geq 0.5\}$ corresponds to the median survival time under a given treatment regime.

In this paper, we study the asymptotic distributions of the estimated value function under the optimal treatment regimes derived. The asymptotic properties of $\hat{\eta}$ in the treatment regime function are very challenging to obtain. The rate of convergence of $\hat{\eta}$ is slower than the classical $n^{1/2}$ -rate due to the indicator function $I(\eta^T \tilde{X} \geq 0)$, and the resulting limiting distribution is not standard. Matsouaka *et al.* (2014) studied a special case where the estimated value function depends on a single threshold value and showed that the estimator of the threshold that maximizes the estimated value function has the $n^{1/3}$ -rate. We conjecture that our estimator $\hat{\eta}$ should also have $n^{1/3}$ -rate. **This is an interesting problem that warrants future research.**

Acknowledgements

The authors are grateful to two referees and the Associate Editor for their thoughtful and suggestive comments, which have helped to improve greatly on an earlier manuscript. The work was partially supported by National Institutes of Health grants R01 CA140632 and P01 CA142538.

Appendix A: Proof of theorems

To establish the asymptotic results that are given in theorems 1 and 2, we need to assume some regularity conditions. Recall that a working logistic model (3) is assumed for the propensity scores with parameters $\boldsymbol{\theta}$ for the IPSWKME and a working proportional hazards model (5) is further assumed for the survival time T for the AIPSWKME with parameters $\boldsymbol{\beta}$ and Λ_0 . Let $\boldsymbol{\nu}_{Ai} = (\mathbf{X}_i^T, A_i, A_i \mathbf{X}_i^T)^T$ and $\boldsymbol{\nu}_{\eta i} = (\mathbf{X}_i^T, g_{\eta}(\mathbf{X}_i), g_{\eta}(\mathbf{X}_i) \mathbf{X}_i^T)^T$. Define

$$K_1^1(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta}) = \int_0^t \frac{(2A - 1) dN(u)}{\pi^* E\{w_{\eta}^* Y(u)\}},$$

$$K_2^1(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta}) = \int_0^t \frac{(2A - 1) Y(u) E[\{(2A - 1) g_{\eta}(\mathbf{X}) + (1 - A)\} dN(u)]}{[\pi^* E\{w_{\eta}^* Y(u)\}]^2},$$

where $w_{\eta}^* = [A g_{\eta}(\mathbf{X}) + (1 - A)\{1 - g_{\eta}(\mathbf{X})\}]/\pi^*$ and $\pi^* = \pi(\mathbf{X}; \boldsymbol{\theta}^*)A + \{1 - \pi(\mathbf{X}; \boldsymbol{\theta}^*)\}(1 - A)$. In addition, define

$$K_1^A(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta}) = \int_0^t \frac{J_1^A(u) - J_0^A(u)}{E[\{L_1^A(u) - L_0^A(u)\}g_{\boldsymbol{\eta}}(\mathbf{X}) + L_0^A(u)]},$$

$$K_2^A(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta}) = \int_0^t \frac{\{L_1^A(u) - L_0^A(u)\}E[\{J_1^A(u) - J_0^A(u)\}g_{\boldsymbol{\eta}}(\mathbf{X}) + J_0^A(u)]}{(E[\{L_1^A(u) - L_0^A(u)\}g_{\boldsymbol{\eta}}(\mathbf{X}) + L_0^A(u)])^2},$$

where

$$J_k^A(u) = \frac{1 - k - (-1)^k A}{\pi^*} dN(u) + e_k \left\{ 1 - \frac{1 - k - (-1)^k A}{\pi^*} \right\} \exp\{-\Lambda_0^*(u)e_k\} S_C(u) d\Lambda_0^*(u),$$

$$L_k^A(u) = \frac{1 - k - (-1)^k A}{\pi^*} Y(u) + \left\{ 1 - \frac{1 - k - (-1)^k A}{\pi^*} \right\} \exp\{-\Lambda_0^*(u)e_k\} S_C(u)$$

and $e_k = \exp\{\boldsymbol{\beta}^{*T}(\mathbf{X}^T, k, k\mathbf{X}^T)^T\}$, $k = 0, 1$. We assume the following conditions.

Assumption 1. The covariates \mathbf{X} are bounded.

Assumption 2. The propensity score $\pi(\mathbf{X})$ is bounded away from 0 and 1 for all possible values of \mathbf{X} .

Assumption 3. The equation

$$E\left[\left\{A - \frac{\exp(\boldsymbol{\theta}^T \tilde{\mathbf{X}})}{1 + \exp(\boldsymbol{\theta}^T \tilde{\mathbf{X}})}\right\} \tilde{\mathbf{X}}\right] = 0$$

has a unique solution $\boldsymbol{\theta}^*$.

Assumption 4. The equation

$$E\left(\int_0^\tau \left[\boldsymbol{\nu}_{Ai} - \frac{E\{Y_i(s) \exp(\boldsymbol{\beta}^T \boldsymbol{\nu}_{Ai}) \boldsymbol{\nu}_{Ai}\}}{E\{Y_i(s) \exp(\boldsymbol{\beta}^T \boldsymbol{\nu}_{Ai})\}}\right] dN_i(s)\right) = 0.$$

has a unique solution $\boldsymbol{\beta}^*$, where $\tau > t$ is a prespecified time point satisfying $P(\tilde{T}_i \geq \tau) > 0$. Let $\Lambda_0^*(u) = E[\int_0^u dN_i(s) / E\{Y_i(s) \exp(\boldsymbol{\beta}^{*T} \boldsymbol{\nu}_{Ai})\}]$ and it satisfies $\Lambda_0^*(\tau) < \infty$.

Assumption 5. $\sup_{\|\boldsymbol{\eta}\|=1} E[\{K_j^1(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta})\}^2] < \infty$ and $\sup_{\|\boldsymbol{\eta}\|=1} E[\{K_j^A(\mathbf{X}, A, \tilde{T}, \delta; \boldsymbol{\eta})\}^2] < \infty$, $j = 1, 2$.

Assumption 6. $nh \rightarrow \infty$ and $nh^4 \rightarrow 0$ as $n \rightarrow \infty$.

Under assumed regularity conditions 1–4, we have the following asymptotic representations:

$$\begin{aligned} \sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}^*) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_{1i} + o_p(1), \\ \sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_{2i} + o_p(1), \\ \sqrt{n}\{\hat{\Lambda}_0(u) - \Lambda_0^*(u)\} &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_{3i}(u) + o_p(1), \\ \sqrt{n}\{\hat{S}_C(u) - S_C(u)\} &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \phi_{4i}(u) + o_p(1), \end{aligned}$$

where the ϕ_{1i} s and ϕ_{2i} s are independently and identically distributed mean 0 vectors, and $\phi_{3i}(u)$ and $\phi_{4i}(u)$ are independent mean 0 processes. Moreover, consistent estimators $\hat{\phi}_{1i}$, $\hat{\phi}_{2i}$, $\hat{\phi}_{3i}(u)$ and $\hat{\phi}_{4i}(u)$ of ϕ_{1i} , ϕ_{2i} , $\phi_{3i}(u)$ and $\phi_{4i}(u)$ can be easily obtained.

In what follows, we give a sketch for the proof of theorem 1. The detailed proofs for theorems 1 and 2 are provided in the on-line supplementary appendix.

A.1. Proof of theorem 1

For any given regime g_η , we first derive the asymptotic properties for the corresponding inverse propensity score weighted Nelson–Aalen estimator. Specifically,

$$\hat{\Lambda}_1(u; \eta) \equiv \hat{\Lambda}_1(u; \eta, \hat{\theta}) = \int_0^u \frac{\sum_{i=1}^n \hat{w}_{\eta_i} dN_i(s)}{\sum_{i=1}^n \hat{w}_{\eta_i} Y_i(s)}. \quad (8)$$

It is easy to show that $\hat{S}_1(u; \eta)$ and $\exp\{-\hat{\Lambda}_1(u; \eta)\}$ are asymptotically equivalent for any given η . Therefore, the asymptotic properties of $\hat{S}_1(u; \eta)$ easily follow those of $\hat{\Lambda}_1(u; \eta)$.

When the propensity score model is correctly specified, we have $\theta^* = \theta$ and $w_{\eta_i}^* = w_{\eta_i}$. Then

$$n^{-1} \sum_{i=1}^n \hat{w}_{\eta_i} Y_i(s) \rightarrow_p E\{w_{\eta_i} Y_i(s)\} = E[Y^*\{g_\eta(\mathbf{X}); s\}]$$

uniformly for $s \in [0, \tau]$ as $n \rightarrow \infty$. Similarly, we have

$$n^{-1} \sum_{i=1}^n \hat{w}_{\eta_i} dN_i(s) \rightarrow_p E\{w_{\eta_i} dN_i(s)\} = E[dN^*\{g_\eta(\mathbf{X}); s\}]$$

uniformly for $s \in [0, \tau]$ as $n \rightarrow \infty$. Therefore,

$$\begin{aligned} \hat{\Lambda}_1(u; \eta) &\rightarrow_p \int_0^u \frac{E[dN^*\{g_\eta(\mathbf{X}); s\}]}{E[Y^*\{g_\eta(\mathbf{X}); s\}]} = \int_0^u \frac{S_C(s) dP[T^*\{g_\eta(\mathbf{X})\} \leq s]}{S_C(s) P[T^*\{g_\eta(\mathbf{X})\} \geq s]} \\ &= -\log\{S^*(u; \eta)\} \equiv \Lambda^*(u; \eta), \end{aligned}$$

which establish the consistency given in part (a) of theorem 1.

Next, we derive the asymptotic distribution of $\hat{\Lambda}_1(u; \eta)$. By applying the first-order Taylor series expansion of $\hat{\Lambda}_1(u; \eta)$ with respect to parameter θ and some empirical process approximation techniques, we have

$$\begin{aligned} \sqrt{n}\{\hat{\Lambda}_1(u; \eta) - \Lambda^*(u; \eta)\} &= n^{-1/2} \sum_{i=1}^n \left(\int_0^u \frac{w_{\eta_i} dM_i^*\{g_\eta(\mathbf{X}); s\}}{E[Y^*\{g_\eta(\mathbf{X}); s\}]} + D_1(u)^T \phi_{1i} \right) + o_p(1) \\ &\equiv n^{-1/2} \sum_{i=1}^n \zeta_i(u; \eta) + o_p(1), \end{aligned}$$

where $M_i^*\{g_\eta(\mathbf{X}); s\} = N_i^*\{g_\eta(\mathbf{X}); s\} - \int_0^s Y_i^*\{g_\eta(\mathbf{X}); v\} d\Lambda^*(v; \eta)$ is a mean 0 martingale process and $D_1(u) = \lim_{n \rightarrow \infty} \partial \hat{\Lambda}_1(u; \eta, \theta)/\partial \theta$. By the delta method, we have

$$\sqrt{n}\{\hat{S}_1(u; \eta) - S^*(u; \eta)\} = -S^*(u; \eta) n^{-1/2} \sum_{i=1}^n \zeta_i(u; \eta) + o_p(1),$$

which converges weakly to a mean 0 Gaussian process by applying empirical process theory. This proves part (b) of theorem 1.

Since $\hat{\eta}_1^{\text{opt}}$ is the maximizer of $\hat{S}_1(t; \eta)$ and η^{opt} is the maximizer of $S^*(t; \eta)$, following similar arguments in Zhang *et al.* (2012b), we have

$$\sqrt{n}\{\hat{S}_1(t; \hat{\eta}_1^{\text{opt}}) - S^*(t; \eta^{\text{opt}})\} - \sqrt{n}\{\hat{S}_1(t; \eta^{\text{opt}}) - S^*(t; \eta^{\text{opt}})\} = o_p(1).$$

It follows that $\sqrt{n}\{\hat{S}_1(t; \hat{\eta}_1^{\text{opt}}) - S^*(t; \eta^{\text{opt}})\} \rightarrow^d N\{0, \Sigma_1(t; \eta^{\text{opt}})\}$, where $\Sigma_1(t; \eta^{\text{opt}}) = S^*(t; \eta^{\text{opt}})^2 E\{\zeta_i^2(t; \eta^{\text{opt}})\}$. This proves part (c) of theorem 1.

Finally, we show that $\hat{S}_1(t; \hat{\eta}_1^{\text{opt}})$ and $\tilde{S}_1(t; \hat{\eta}_1^{\text{opt}})$ are asymptotically equivalent. For any given η , we have

$$\sqrt{n}\{\tilde{\Lambda}_1(t; \eta) - \hat{\Lambda}_1(t; \eta)\} = \sqrt{n} \frac{1}{n} \sum_{i=1}^n \left\{ \Phi\left(\frac{\eta^T \mathbf{X}_i}{h}\right) - I(\eta^T \mathbf{X}_i \geq 0) \right\} \times K_1^1(\mathbf{X}_i, A_i, \tilde{T}_i, \delta; \eta) \quad (9)$$

$$\begin{aligned} &+ \sqrt{n} \frac{1}{n} \sum_{i=1}^n \left\{ \Phi\left(\frac{\eta^T \mathbf{X}_i}{h}\right) - I(\eta^T \mathbf{X}_i \geq 0) \right\} \times K_2^1(\mathbf{X}_i, A_i, \tilde{T}_i, \delta; \eta) \quad (10) \\ &+ o_p(1). \end{aligned}$$

Following similar arguments to those in Heller (2007), we can show that $\sup_{\|\eta\|=1} |(9)| = o_p(1)$ and $\sup_{\|\eta\|=1} |(10)| = o_p(1)$. Therefore, we have $\sqrt{n}\{\tilde{\Lambda}_1(t; \eta) - \hat{\Lambda}_1(t; \eta)\} = o_p(1)$ uniformly in η , which implies that $\sqrt{n}\{\tilde{S}_1(t; \eta) - \hat{S}_1(t; \eta)\} = o_p(1)$ uniformly in η . It follows that $\sqrt{n}\{\tilde{S}_1(t; \hat{\eta}_1^{\text{opt}}) - \hat{S}_1(t; \hat{\eta}_1^{\text{opt}})\} = o_p(1)$, which proves part (d) of theorem 1.

References

- Bai, X., Tsiatis, A. A., Lu, W. and Song, R. (2014) Optimal treatment regimes for survival endpoints using a locally-efficient doubly-robust estimator from a classification perspective. *Technical Report*. Department of Statistics, North Carolina State University, Raleigh.
- Cheng, S. C., Wei, L. J. and Ying, Z. (1995) Analysis of transformation models with censored data. *Biometrika*, **82**, 835–845.
- Cox, D. R. (1972) Regression models and life-tables (with discussion). *J. R. Statist. Soc. B*, **34**, 187–220.
- Goldberg, Y. and Kosorok, M. R. (2012) Q-learning with censored data. *Ann. Statist.*, **40**, 529–560.
- Hammer, S. M., Katzenstein, D. A., Hughes, M. D., Gundacker, H., Schooley, R. T., Haubrich, R. H., Henry, W. K., Lederman, M. M., Phair, J. P., Niu, M., Hirsch, M. S. and Merigan, T. C. (1996) A trial comparing nucleoside monotherapy with combination therapy in hiv-infected adults with cd4 cell counts from 200 to 500 per cubic millimeter. *New Engl. J. Med.*, **335**, 1081–1090.
- Heller, G. (2007) Smoothed rank regression with censored data. *J. Am. Statist. Ass.*, **102**, 552–559.
- Matsouaka, R. A., Li, J. and Cai, T. (2014) Evaluating marker-guided treatment selection strategies. *Biometrics*, **70**, 489–499.
- Mebane, Jr, W. R. and Sekhon, J. S. (2011) Genetic optimization using derivatives: the rgenoud package for R. *J. Statist. Softw.*, **42**, 1–26.
- Murphy, S. A. (2003) Optimal dynamic treatment regimes. *J. R. Statist. Soc. B*, **65**, 331–355.
- Murphy, S. A. (2005) An experimental design for the development of adaptive treatment strategies. *Statist. Med.*, **24**, 1455–1481.
- Robins, J. M. (2004) Optimal structural nested models for optimal sequential decisions. In *Proc. 2nd Seattle Symp. Biostatistics* (eds D. Y. Lin and P. J. Heagerty), pp. 189–326. New York: Springer.
- Rubin, D. B. (1974) Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.*, **66**, 688–701.
- Watkins, C. and Dayan, P. (1992) Q-learning. *Mach. Learn.*, **8**, 279–292.
- Watkins, C. J. (1989) Learning from delayed rewards. *PhD Thesis*. University of Cambridge, Cambridge.
- Zhang, B., Tsiatis, A. A., Davidian, M., Zhang, M. and Laber, E. B. (2012a) Estimating optimal treatment regimes from a classification perspective. *Stat.*, **1**, 103–114.
- Zhang, B., Tsiatis, A. A., Laber, E. B. and Davidian, M. (2012b) A robust method for estimating optimal treatment regimes. *Biometrics*, **68**, 1010–1018.
- Zhang, B., Tsiatis, A. A., Laber, E. B. and Davidian, M. (2013) Robust estimation of optimal dynamic treatment regimes for sequential treatment decisions. *Biometrika*, **100**, 681–694.
- Zhao, Y., Kosorok, M. R. and Zeng, D. (2009) Reinforcement learning design for cancer clinical trials. *Statist. Med.*, **28**, 3294–3315.
- Zhao, L., Tian, L., Cai, T., Claggett, B. and Wei, L. J. (2013) Effectively selecting a target population for a future comparative study. *J. Am. Statist. Ass.*, **108**, 527–539.
- Zhao, Y., Zeng, D., Laber, E., Song, R., Yuan, M. and Kosorok, M. (2015) Doubly robust learning for estimating individualized treatment with censored data. *Biometrika*, **102**, 151–168.
- Zhao, Y., Zeng, D., Rush, A. J. and Kosorok, M. R. (2012) Estimating individualized treatment rules using outcome weighted learning. *J. Am. Statist. Ass.*, **107**, 1106–1118.

Supporting information

Additional ‘supporting information’ may be found in the on-line version of this article:

‘Supplement appendix for “On estimation of optimal treatment regimes for maximizing t -year survival probability”’.