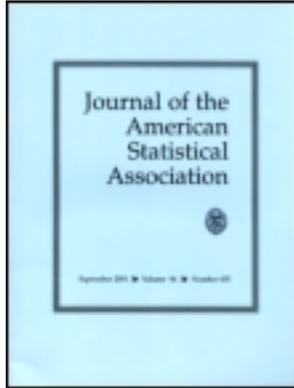


This article was downloaded by: [Texas A&M University Libraries]

On: 08 April 2013, At: 15:03

Publisher: Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK



Journal of the American Statistical Association

Publication details, including instructions for authors and subscription information:

<http://amstat.tandfonline.com/loi/uasa20>

Nonparametric Estimation for Censored Mixture Data With Application to the Cooperative Huntington's Observational Research Trial

Yuanjia Wang^a, Tanya P. Garcia^b & Yanyuan Ma^c

^a Department of Biostatistics, Columbia University, New York, NY, 10032

^b Department of Statistics, Texas A&M University, College Station, TX, 77843-314

^c Department of Statistics, Texas A&M University, College Station, TX, 77843-3143

Version of record first published: 21 Dec 2012.

To cite this article: Yuanjia Wang, Tanya P. Garcia & Yanyuan Ma (2012): Nonparametric Estimation for Censored Mixture Data With Application to the Cooperative Huntington's Observational Research Trial, *Journal of the American Statistical Association*, 107:500, 1324-1338

To link to this article: <http://dx.doi.org/10.1080/01621459.2012.699353>

PLEASE SCROLL DOWN FOR ARTICLE

Full terms and conditions of use: <http://amstat.tandfonline.com/page/terms-and-conditions>

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden.

The publisher does not give any warranty express or implied or make any representation that the contents will be complete or accurate or up to date. The accuracy of any instructions, formulae, and drug doses should be independently verified with primary sources. The publisher shall not be liable for any loss, actions, claims, proceedings, demand, or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of this material.

Nonparametric Estimation for Censored Mixture Data With Application to the Cooperative Huntington's Observational Research Trial

Yuanjia WANG, Tanya P. GARCIA, and Yanyuan MA

This work presents methods for estimating genotype-specific outcome distributions from genetic epidemiology studies where the event times are subject to right censoring, the genotypes are not directly observed, and the data arise from a mixture of scientifically meaningful subpopulations. Examples of such studies include kin-cohort studies and quantitative trait locus (QTL) studies. Current methods for analyzing censored mixture data include two types of nonparametric maximum likelihood estimators (NPMLEs; Type I and Type II) that do not make parametric assumptions on the genotype-specific density functions. Although both NPMLEs are commonly used, we show that one is inefficient and the other inconsistent. To overcome these deficiencies, we propose three classes of consistent nonparametric estimators that do not assume parametric density models and are easy to implement. They are based on inverse probability weighting (IPW), augmented IPW (AIPW), and nonparametric imputation (IMP). AIPW achieves the efficiency bound without additional modeling assumptions. Extensive simulation experiments demonstrate satisfactory performance of these estimators even when the data are heavily censored. We apply these estimators to the Cooperative Huntington's Observational Research Trial (COHORT), and provide age-specific estimates of the effect of mutation in the Huntington gene on mortality using a sample of family members. The close approximation of the estimated noncarrier survival rates to that of the U.S. population indicates small ascertainment bias in the COHORT family sample. Our analyses underscore an elevated risk of death in Huntington gene mutation carriers compared with that in noncarriers for a wide age range, and suggest that the mutation equally affects survival rates in both genders. The estimated survival rates are useful in genetic counseling for providing guidelines on interpreting the risk of death associated with a positive genetic test, and in helping future subjects at risk to make informed decisions on whether to undergo genetic mutation testing. Technical details and additional numerical results are provided in the online supplementary materials.

KEY WORDS: Censored data; Finite mixture model; Huntington's disease; Kin-cohort design; Quantitative trait locus.

1. INTRODUCTION

In some genetic epidemiology studies, a research goal is to estimate genotype-specific cumulative distributions of an outcome from mixture data of scientifically meaningful subpopulations where genotypes are not directly observed. Examples of such studies include kin-cohort studies (Struewing et al. 1997; Wacholder et al. 1998; Wang et al. 2008; Mai et al. 2009) and quantitative trait locus (QTL) studies (Lander and Botstein 1989; Wu, Ma, and Casella 2007). In kin-cohort studies, scientists sample and genotype an initial cohort of subjects (proband), possibly enriched with mutation carriers. They then collect family history of the disease (phenotype) from family members of the probands through systematic and validated interviews of the probands (Marder et al. 2003). While it is impractical and costly to interview family members in-person to collect their blood samples and obtain genotypes, it is possible to calculate the probability of each relative having a certain genotype based

on the relationship with the proband and the proband's genotype. Thus, kin-cohort studies differ from other types of case-control family studies (Li, Yang, and Schwartz 1998) in that genetic information in family members is not readily available. Distributions of the observed phenotypes in the relatives are therefore a mixture of genotype-specific distributions.

In the interval mapping of quantitative traits (Lander and Botstein 1989), the genotype of a QTL is not observed, so trait distributions are mixtures of the QTL genotype-specific distributions. The mixing proportions are computed based on the observed flanking marker genotypes and recombination fractions between the marker and the putative QTL. In many controlled QTL experiments such as backcross or intercross, the mixing proportions can be easily obtained, and interest is in estimating the genotype-specific distributions.

The unobserved genotype information in both kin-cohort and QTL studies makes inference of genotype-specific distributions difficult. Inference is further complicated by right censoring as patients in the study may drop out or become lost to follow-up. The focus of the current article is to develop simple, robust, and efficient estimators to improve upon the available methods in the literature for analyzing such censored mixture data.

Many statistical methods have been developed for modeling and analyzing censored mixture data in QTL mappings and kin-cohort studies. Sometimes, the biological underpinning of the development of a disease trait suggests a suitable parametric

Yuanjia Wang is Assistant Professor, Department of Biostatistics, Columbia University, New York, NY 10032 (E-mail: yuanjia.wang@columbia.edu). Tanya P. Garcia is Post-Doctoral Fellow, Department of Statistics, Texas A&M University, College Station, TX 77843-314 (E-mail: tpgarcia@stat.tamu.edu). Yanyuan Ma is Professor, Department of Statistics, Texas A&M University, College Station, TX 77843-3143 (E-mail: ma@stat.tamu.edu). This research is supported in part by the National GEM (Graduate Degrees for Minorities in Engineering and Science) Consortium, the Philanthropic Education Organization (PEO) Scholarship Award, and grants from the U.S. National Science Foundation (0906341 and 1206693), the National Institutes of Health (NS073671-01 and AG031113-01A2), and the National Cancer Institute (R25T-CA090301). Samples and data from the Cooperative Huntington's Observational Research Trial (COHORT) study, which receives support from HP Therapeutics, Inc., were used in this study. The authors thank the Huntington Study Group COHORT investigators and coordinators who collected data and/or samples used in this study, as well as participants and their families who made this work possible.

function that offers meaningful interpretation of the biological structure (Wu et al. 2000). In these cases, it is reasonable to use maximum likelihood-based parametric methods (Lander and Botstein 1989; Wu et al. 2002). In some QTL experiments, a semiparametric Cox proportional hazards model may also be suitable (Diao and Lin 2005; Zeng and Lin 2007), but a proportional hazards assumption is not always valid, such as in some applications with Huntington's disease (HD) data (Langbehn et al. 2004). In fact, in many situations, there may not be sufficient biological knowledge to warrant particular parametric or semiparametric models; hence, concerns of model misspecification naturally arise. To alleviate these issues, more flexible nonparametric estimation of the distribution functions becomes essential (Yu and Lin 2008; Zhao and Wu 2008). Throughout this work, the term "nonparametric" refers to leaving the probability density (or hazard) functions completely unspecified.

For QTL data, Fine, Zhou, and Yandell (2004) developed a nonparametric method that exploits the property of independence between the censoring and the event of interest. Wang et al. (2007) proposed a nonparametric method for kin-cohort data when the censoring times are observed for all subjects. When censoring times are random and are not observed for all subjects, Wacholder et al. (1998) proposed a nonparametric maximum likelihood estimator (Type I NPMLE) consisting of a combination of several NPMLEs and a linear transformation. Chatterjee and Wacholder (2001) proposed a direct maximization of the nonparametric likelihood (Type II NPMLE) with respect to the conditional distributions and used an expectation-maximization (EM) algorithm to find the maximizer. Although in many situations, NPMLEs are consistent and even efficient, we demonstrate the surprising result that the Type I is highly inefficient and the Type II is inconsistent.

To overcome the shortcomings of the aforementioned methods, we provide several consistent and efficient nonparametric estimators by casting this problem in a missing-data framework. Given a complete-data influence function when there is no censoring (see the Appendix; Ma and Wang 2012), we propose an inverse probability weighting (IPW) estimator and derive an optimal augmentation term to obtain the optimal estimator. We demonstrate that the optimal augmented IPW (AIPW) estimator achieves the efficiency bound without extra modeling assumptions or complicated computational procedures. We also propose an imputation (IMP) estimator that is easy to implement and does not require additional modeling assumptions for the imputation step.

The rest of the article is organized as follows. Section 1.1 presents a large collaborative study of HD to which we apply our proposed estimators. Section 2 describes the inefficiency and inconsistency of the two existing NPMLE methods. To improve upon these methods, we propose several nonparametric estimators in Section 3 that are consistent, efficient, and easy to implement. We demonstrate the asymptotic properties of these estimators and examine their finite-sample performance through comprehensive simulation studies in Section 4. The methods are applied to the HD study in Section 5, and Section 6 concludes the article with some discussion. The technical details and additional numerical results are provided in the Appendix and in the online supplementary materials, with tables

and figures in the supplementary materials indicated with a preceding "S."

1.1 The Cooperative Huntington's Observational Research Trial (COHORT)

Huntington's disease is a degenerative, genetic disorder that targets nerve cells in the brain and leads to cognitive decline, involuntary muscle spasms, and psychological problems. Affected individuals typically begin to see neurological and physical symptoms around 30–50 years of age, and eventually die from pneumonia, heart failure, or other complications 15–20 years after the disease onset (Foroud et al. 1999). The severity of the disease has prompted the development of several organizations, such as the Huntington Study Group (<http://www.huntington-study-group.org/>), which are devoted to studying the causes, effects, and possible treatments for HD. A particular study organized by roughly 42 Huntington Study Group research centers in North America and Australia is the Cooperative Huntington's Observational Research Trial (COHORT; Dorsey et al. 2008). Since 2005, the principal investigators of COHORT have been collecting ongoing information from affected or at-risk adults and their family members 15 years of age and older.

Huntington's disease is caused by unstable CAG repeats expansion in the Huntington gene (Huntington's Disease Collaborative Research Group 1993). In a genetic counseling setting, CAG repeats ≥ 36 is defined as positive for HD gene mutation, or carrier, and CAG < 36 is defined as negative, or noncarrier (Rubinsztein et al. 1996). Proband participants in COHORT undergo a clinical evaluation where blood samples are genotyped for being a carrier or noncarrier of HD mutation. While the HD mutation status is ascertained in probands, high costs of in-person interviews of family members prevents collection of their blood samples. Family members' morbidity and mortality information, such as age at death, is obtained through a systematic interview of the probands. Although a relative's HD gene mutation status is unavailable, the probability of carrying a mutation can still be obtained based on the relative's relationship with the proband and the proband's mutation status (Khoury, Beaty, and Cohen 1993, sec. 8.4). The distribution of the relative's age at death is therefore a mixture of the genotype-specific distributions with known, subject-specific mixing proportions.

Despite the identification of the causative gene, there is currently no effective treatment that modifies HD progression. One of the goals of COHORT is to estimate the risks of adverse events, such as disease onset or death, associated with carrying a mutation, and to use these parameters to design future clinical trials for intervention or treatment of HD. For example, the power calculation of a trial with survival as the primary endpoint will depend on parameters such as risk ratio in carriers and noncarriers. The proposed methods here can aid in estimating these important parameters, and also has benefits in genetic counseling for patients and their family members. The estimated survival function in HD mutation carriers provides guidelines on interpreting the risk of death associated with a positive genetic mutation test, and facilitates subjects at risk to make important life decisions, such as marriage or having children. We show some examples of the utilities of the survival estimates in Section 5.

2. SOME EXISTING NONPARAMETRIC ESTIMATORS FOR CENSORED MIXTURE DATA

We consider censored mixture data denoted as triplets $(\mathbf{Q}_i = \mathbf{q}_i, X_i = x_i, \Delta_i = \delta_i)$, which are independent and identically distributed. For the i th subject, \mathbf{Q}_i is a p -dimensional vector of the random mixture proportions computed from available genotype data on the proband in kin-cohort studies or from flanking markers in QTL studies. In a kin-cohort study, \mathbf{Q}_i may be a two-dimensional vector, where Q_{i1} represents the probability of being a mutation carrier, and Q_{i2} a noncarrier. To illustrate the computation of \mathbf{Q}_i , let L_i denote the unobserved genotype in a relative, and let L_{i0} denote the observed genotype in a proband. Let p_A denote the population frequency of the mutation allele A , and let a denote the wild type. Consider a heterozygous carrier proband with genotype Aa . Assuming Mendelian transmission, the probability of a parent of a proband being a carrier is $Q_{i1} = \Pr(L_i = AA \text{ or } Aa | L_{i0} = Aa) = \frac{1}{2}(1 + p_A)$. The probability for a sibling of a proband to be a carrier is $Q_{i1} = \Pr(L_i = AA \text{ or } Aa | L_{i0} = Aa) = -\frac{1}{4}p_A^2 + \frac{3}{4}p_A + \frac{1}{2}$. When $p_A \approx 0$, the two Q_{i1} 's are both $\frac{1}{2}$. The \mathbf{Q}_i for other types of relatives and other types of probands (homozygous or non-carrier probands) are computed similarly; see Khoury, Beaty, and Cohen (1993, sec. 8.4) for details.

In general, for both QTL and kin-cohort studies, \mathbf{Q}_i has a discrete distribution with a finite support, say $\mathbf{u}_1, \dots, \mathbf{u}_m$. Its probability mass function, denoted as $p_{\mathbf{Q}}$, is determined by the experimental design. For example, in a backcross QTL study, \mathbf{Q}_i is a two-dimensional vector that will take four possible values $(1, 0)^T, (\theta, 1 - \theta)^T, (1 - \theta, \theta)^T$, and $(0, 1)^T$, where θ is the known recombination fraction between the putative QTL and the flanking marker. The probability of \mathbf{Q}_i taking these four values is determined by the marker genotype frequencies computed from the observed marker data (e.g., Wu, Ma, and Casella 2007, table 10.4). In kin-cohort studies, the distribution of \mathbf{Q}_i is determined by the type of relatives collected (e.g., parents, siblings, and children) and the distribution of the probands' genotypes (e.g., number of noncarrier probands, heterozygote probands, and homozygous probands).

Last, $X_i = \min(T_i, C_i)$, where T_i is a subject's event time and C_i is a random continuous censoring time independent of T_i ; and $\Delta_i = I(T_i \leq C_i)$ is the censoring indicator. We let $f(\cdot)$ denote the p -dimensional unspecified conditional probability density function of T , given genotypes in p genotype groups, and let $F(\cdot)$ denote the corresponding cumulative distribution function. Interest lies in estimating $F(t)$ for any fixed time t . In the COHORT study, we have $p = 2$, with $F_1(t)$ and $F_2(t)$ corresponding to the age-at-death distribution for HD gene mutation carriers and noncarriers, respectively. Throughout, except when specifically pointed out, we assume that the event times x_1, \dots, x_n have no ties, and that the censoring distribution is common for all subjects. Then, letting $G(\cdot)$ denote the survival function of C and $g(\cdot)$ its corresponding density, the log-likelihood of n observations is

$$\sum_{i=1}^n \log(p_{\mathbf{Q}}(\mathbf{q}_i) \{ \mathbf{q}_i^T \mathbf{f}(x_i) G(x_i) \}^{\delta_i} [\{ 1 - \mathbf{q}_i^T \mathbf{F}(x_i) \} \mathbf{g}(x_i)]^{1-\delta_i}), \tag{1}$$

where we use the fact that $\mathbf{q}_i^T \mathbf{1}_p = 1$, with $\mathbf{1}_p$ being a p -dimensional vector of 1's.

2.1 The Type I NPMLE and Its Inefficiency

The Type I NPMLE was proposed in the literature to analyze kin-cohort data (Wacholder et al. 1998). It first maximizes (1) with respect to $\mathbf{q}_i^T \mathbf{f}(x_i)$'s, then recovers $F(t)$ through a linear transformation. Although an NPMLE-based estimator is usually efficient, it is not so for mixture data, and the magnitude of efficiency loss can be large.

To describe the Type I NPMLE, we reformulate the maximization problem by evoking the assumption that \mathbf{Q} has finite support $\mathbf{u}_1, \dots, \mathbf{u}_m$ and by letting $s_j(x_k) = \mathbf{u}_j^T \mathbf{f}(x_k)$ and $S_j(x_k) = 1 - \mathbf{u}_j^T \mathbf{F}(x_k)$. Throughout this work, we refer to the p different genotype populations as p subpopulations, and the m different \mathbf{u}_j values as m subgroups. In the literature, the Type I NPMLE assumes random censoring, and hence, the censoring distribution does not contribute information to the parameter of interest. Therefore, ignoring $G(\cdot)$ and $g(\cdot)$ in Equation (1), it maximizes the equivalent target function

$$\sum_{j=1}^m \sum_{i=1}^n \log \{ s_j(x_i)^{\delta_i} S_j(x_i)^{1-\delta_i} \} I(\mathbf{q}_i = \mathbf{u}_j) \tag{2}$$

with respect to $s_j(x_i)$'s and subject to $\sum_{i=1}^n s_j(x_i) I(\mathbf{q}_i = \mathbf{u}_j) \leq 1, s_j(x_i) \geq 0$ for $j = 1, \dots, m$. This is equivalent to m separate maximization problems, each concerning $s_j(\cdot)$ and $S_j(\cdot)$ only, so the maximizers are the classical Kaplan–Meier estimators. That is,

$$\widehat{S}_j(a) = \prod_{x_i \leq a, \mathbf{q}_i = \mathbf{u}_j} \left\{ 1 - \frac{\delta_i}{\sum_{\mathbf{q}_k = \mathbf{u}_j} I(x_k \geq x_i)} \right\}$$

and $s_j(a) = S_j(a^-) - S_j(a)$ for all a . Using the linear relation $\mathbf{u}_j^T \mathbf{F}(t) = 1 - S_j(t)$ for $j = 1, \dots, m$, we then recover the Type I NPMLE estimator as

$$\widetilde{F}(t) = (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T \{ \mathbf{1}_m - \widehat{\mathbf{S}}(t) \},$$

where $\widehat{\mathbf{S}}(t) = \{ \widehat{S}_1(t), \dots, \widehat{S}_m(t) \}^T$ and $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_m)^T$. In this notation, $\mathbf{S}(t) = \mathbf{1}_m - \mathbf{U} \mathbf{F}(t)$. The consistency of the Kaplan–Meier estimator of $\mathbf{S}(t)$ ensures the consistency of $\widetilde{F}(t)$. The inefficiency of $\widetilde{F}(t)$, however, is evident, considering that $\widetilde{F}_w(t) = (\mathbf{U}^T \boldsymbol{\Sigma}^{-1} \mathbf{U})^{-1} \mathbf{U}^T \boldsymbol{\Sigma}^{-1} \{ \mathbf{1}_m - \widehat{\mathbf{S}}(t) \}$, with $\boldsymbol{\Sigma}$ denoting the variance-covariance matrix of $\widehat{\mathbf{S}}(t)$, yields a more efficient estimator. In this case, $\boldsymbol{\Sigma}$ is a diagonal matrix because each of the m components of $\widehat{\mathbf{S}}(t)$ is estimated using a distinct subset of the observations. Hence, $\widetilde{F}_w(t)$ is a weighted version of the Type I NPMLE, and this simple weighting scheme improves the estimation efficiency.

2.2 The Type II NPMLE and Its Inconsistency

The Type II NPMLE is considered an improvement over the Type I NPMLE (Chatterjee and Wacholder 2001). It maximizes the same log-likelihood (1), but with respect to $f(x_i)$'s subject to $\sum_{i=1}^n \mathbf{f}(x_i) \leq \mathbf{1}_p$ and $\mathbf{f}(x_i) \geq \mathbf{0}$ componentwise. Like the Type I NPMLE, the Type II NPMLE also assumes random censoring and ignores $G(\cdot)$ and $g(\cdot)$ in Equation (1). In general, no closed-form solution exists, and the EM algorithm is implemented to obtain the $F(x_i)$'s. Specifically,

we regard the genotypes $L_i = 1, \dots, p$ as missing data and derive the complete-data log-likelihood of the observations $\mathbf{o}_i = (L_i = l_i, X_i = x_i, \Delta_i = \delta_i), i = 1, \dots, n$, as

$$\begin{aligned} \mathcal{L}_{\text{TypeII}}^{\text{comp}}\{\mathbf{o}_1, \dots, \mathbf{o}_n; \mathbf{f}(x_i), \mathbf{F}(x_i), i = 1, \dots, n\} \\ = \sum_{i=1}^n \sum_{k=1}^p I(l_i = k) \log [f_k(x_i)^{\delta_i} \{1 - F_k(x_i)\}^{1-\delta_i}]. \end{aligned}$$

The EM algorithm is an iterative procedure. At the b th iteration, we take the conditional expectation of the complete-data log-likelihood given the observed data (e.g., $\{(X_i = x_i, \Delta_i = \delta_i), i = 1, \dots, n\}$), and update the E-step via

$$\begin{aligned} E[\mathcal{L}_{\text{TypeII}}^{\text{comp}}\{\mathbf{o}_1, \dots, \mathbf{o}_n; \mathbf{f}(x_i), \mathbf{F}(x_i), i = 1, \dots, n\} \\ | \mathbf{f}^{(b)}(x_i), \mathbf{F}^{(b)}(x_i), i = 1, \dots, n] \\ = \sum_{k=1}^p \sum_{i=1}^n \left[\frac{\delta_i q_{ik} f_k^{(b)}(x_i)}{\sum_{k=1}^p q_{ik} f_k^{(b)}(x_i)} \log f_k(x_i) \right. \\ \left. + \frac{(1 - \delta_i) q_{ik} \{1 - F_k^{(b)}(x_i)\}}{\sum_{k=1}^p q_{ik} \{1 - F_k^{(b)}(x_i)\}} \log \{1 - F_k(x_i)\} \right]. \end{aligned}$$

The M-step maximizes the above expression with respect to $\mathbf{f}(x_i)$ and $\mathbf{F}(x_i)$'s subject to $\mathbf{f}(x_i) \geq \mathbf{0}$ and $\mathbf{1} \geq \mathbf{F}(x_i) \geq \mathbf{0}$. To this end, let

$$c_{ik}^{(b)} = \delta_i \frac{q_{ik} f_k^{(b)}(x_i)}{\sum_{k=1}^p q_{ik} f_k^{(b)}(x_i)} + (1 - \delta_i) \frac{q_{ik} \{1 - F_k^{(b)}(x_i)\}}{\sum_{k=1}^p q_{ik} \{1 - F_k^{(b)}(x_i)\}}$$

denote the known quantity based on the b th iteration. Then, the M-step reduces to p separate maximization problems of the form

$$\sum_{i=1}^n c_{ik}^{(b)} [\delta_i \log f_k(x_i) + (1 - \delta_i) \log \{1 - F_k(x_i)\}],$$

for $k = 1, \dots, p$. Viewing this as the log-likelihood of weighted observations, where the i th observation represents $c_{ik}^{(b)}$ observations of the same value, the maximizer is a modified Kaplan–Meier estimator:

$$\begin{aligned} 1 - \check{F}_k^{(b+1)}(t) &= \prod_{x_i \leq t, \delta_i = 1} \left\{ 1 - \frac{\sum_{j=1}^n I(x_j = x_i, \delta_j = 1) c_{jk}^{(b)}}{\sum_{j=1}^n c_{jk}^{(b)} I(x_j \geq x_i)} \right\} \\ &= \prod_{x_i \leq t, \delta_i = 1} \left\{ 1 - \frac{c_{ik}^{(b)}}{\sum_{j=1}^n c_{jk}^{(b)} I(x_j \geq x_i)} \right\}. \end{aligned}$$

Iterating the E- and the M-step until convergence leads to the Type II estimator.

As natural as the Type II NPMLLE appears, we show in Section S.1 of the online supplementary materials the surprising result that it is an inconsistent estimator of $F(t)$.

3. PROPOSED NONPARAMETRIC ESTIMATORS FOR CENSORED MIXTURE DATA

3.1 The IPW and the Optimal AIPW Estimators

To compensate for the deficiencies of the NPMLLEs, we propose a class of nonparametric estimators based on IPW and its augmented version that are consistent and easy to implement. We describe these estimators in terms of their corresponding influence functions.

3.1.1 Inverse Probability Weighting. The notion of IPW was first introduced by Horvitz and Thompson (1952) in the context of survey sampling, and later by Robins, Rotnitzky, and Zhao (1994) in the context of missing data as a means for upweighting subjects who are underrepresented because of missingness. Bang and Tsiatis (2000, 2002) used the IPW to estimate the mean and median medical costs by capturing information from patients whose medical costs were subject to right censoring. In this spirit, we elicit information from the censored observations in the mixture data with an IPW estimator. Specifically, our IPW estimator solves

$$n^{-1} \sum_{i=1}^n \frac{\delta_i \phi(\mathbf{q}_i, x_i)}{\widehat{G}(x_i)} = 0, \tag{3}$$

where ϕ denotes a general influence function for noncensored mixture data corresponding to $\delta_i = 1$ ($i = 1, \dots, n$) in Equation (1) (see the Appendix for elaborations on ϕ) and $\widehat{G}(t)$ is the Kaplan–Meier estimator of $G(t)$:

$$\widehat{G}(t) = \prod_{x_i \leq t} \left\{ 1 - \frac{1 - \delta_i}{\sum_{j=1}^n I(x_j \geq x_i)} \right\}.$$

The intuition behind Equation (3) is that for any subject randomly selected from the population with $X_i = x_i$, the probability that such a subject will not be censored is $G(x_i)$. Therefore, any uncensored subject with $X_i = x_i$ can be regarded as representing $1/G(x_i)$ subjects from the population. By inversely weighting all uncensored subjects with their corresponding probabilities of not being censored, we obtain a consistent estimating equation in (3).

We now characterize the asymptotic behavior of the IPW estimator in terms of its influence function. Let $Y_i(u) = I(X_i \geq u)$, $Y(u) = \sum_{i=1}^n Y_i(u)$, $N_i^c(u) = I(X_i \leq u, \Delta_i = 0)$, and $\lambda^c(\cdot)$ be the hazard function for the censoring distribution. Also, let

$$M_i^c(u) = N_i^c(u) - \int_0^u I(X_i \geq s) \lambda^c(s) ds$$

denote the censoring martingale, and define

$$\mathcal{B}(\mathbf{h}, u) = E\{\mathbf{h}(\cdot) | T_i \geq u\} = \frac{E\{\mathbf{h}(\cdot) I(T_i \geq u)\}}{S(u)},$$

where \mathbf{h} is any p -length function. Then, as derived in Section S.2 of the supplementary materials, the i th influence function for the IPW estimator is

$$\phi_{\text{ipw}}(\mathbf{q}_i, x_i, \delta_i) = \phi(\mathbf{q}_i, t_i) - \int \frac{dM_i^c(u)}{G(u)} \{\phi(\mathbf{q}_i, t_i) - \mathcal{B}(\phi, u)\}.$$

The two terms in ϕ_{ipw} are uncorrelated, given that $\phi(\mathbf{q}_i, x_i)$ is $\mathcal{F}(0)$ measurable, where $\mathcal{F}(u)$ is a filtration defined by the set of σ -algebras generated by $\sigma\{\mathbf{q}_i, I(C_i \leq r), r \leq u; I(T_i \leq x), 0 \leq x < \infty, i = 1, \dots, n\}$. Hence, the estimation variance of the IPW estimator is

$$\begin{aligned} V_{\text{ipw}} &= \text{cov}\{\phi(\mathbf{Q}_i, T_i)\} \\ &+ E \left\{ \int \frac{\mathcal{B}(\phi^{\otimes 2}, u) - \mathcal{B}(\phi, u)^{\otimes 2}}{G^2(u)} \lambda^c(u) Y_i(u) du \right\}, \end{aligned}$$

and a consistent estimator is

$$\widehat{V}_{ipw} = n^{-1} \sum_{i=1}^n \frac{\delta_i \boldsymbol{\phi}(\mathbf{q}_i, x_i) \boldsymbol{\phi}^T(\mathbf{q}_i, x_i)}{\widehat{G}(x_i)} + n^{-1} \sum_{i=1}^n \int \frac{\widehat{B}_1(\boldsymbol{\phi}^{\otimes 2}, u) - \widehat{B}_1(\boldsymbol{\phi}, u)^{\otimes 2}}{\widehat{G}^2(u)} dN_i^c(u),$$

where $\widehat{B}_1(\mathbf{h}, u) = \frac{1}{nS(u)} \sum_{i=1}^n \frac{\delta_i \mathbf{h}(\mathbf{q}_i, x_i, \delta_i) I(x_i \geq u)}{\widehat{G}(x_i)}$ for an arbitrary function $\mathbf{h}(\mathbf{q}_i, x_i, \delta_i)$.

3.1.2 Augmented Inverse Probability Weighting. Although intuitive and easy to implement, the IPW estimator is inefficient. Instead, using a modification motivated by Robins and Rotnitzky (1992), one may adjust the IPW estimator to improve its efficiency. With $\boldsymbol{\phi}$ as the complete-data influence function, Robins and Rotnitzky (1992) provided the following general class of influence functions for censored data:

$$\boldsymbol{\phi}(\mathbf{q}_i, t_i) - \int \frac{dM_i^c(u)}{G(u)} \{\boldsymbol{\phi}(\mathbf{q}_i, t_i) - \mathcal{B}(\boldsymbol{\phi}, u)\} + \int \frac{dM_i^c(u)}{G(u)} [\mathbf{h}(\bar{\mathbf{a}}_i(u), u) - \mathcal{B}(\mathbf{h}, u)]. \quad (4)$$

For our mixture-data problem, $\mathbf{a}_i(u) = \{\mathbf{q}_i, I(u < T_i)\}$ and $\bar{\mathbf{a}}_i(u)$ contains the functions $\mathbf{a}_i(\tilde{u})$ for all $\tilde{u} \leq u$. Compared with the influence function for the IPW estimator, the estimator from (4) contains an augmentation term that may improve the estimation efficiency, and is thus, termed the AIPW estimator. Among all the choices for \mathbf{h} , Robins, Rotnitzky, and Zhou (1994) and van der Laan and Hubbard (1998) showed that

$$\mathbf{h}_{\text{eff}}^* \{\bar{\mathbf{a}}_i(u), u\} = E\{\boldsymbol{\phi}(\mathbf{Q}_i, T_i) | T_i \geq u, \bar{\mathbf{a}}_i(u)\} = \{I(u < X_i) + I(u = X_i, \delta_i = 0)\} E\{\boldsymbol{\phi}(\mathbf{Q}_i, T_i) | \mathbf{q}_i, T_i \geq u\} + I(u = X_i) \delta_i \boldsymbol{\phi}(\mathbf{q}_i, u),$$

with $u \leq X_i$, yields the optimal efficiency. Denoting $\mathbf{h}_{\text{eff},i}(u) = E\{\boldsymbol{\phi}(\mathbf{Q}_i, T_i) | \mathbf{q}_i, T_i \geq u\}$, we have that $\mathbf{h}_{\text{eff}}^* \{\bar{\mathbf{a}}_i(u), u\}$ and $\mathbf{h}_{\text{eff},i}(u)$ are identical except when $u = X_i$ and $\delta_i = 1$. The functional $\mathbf{h}_{\text{eff}}^*$ only appears in the censoring martingale integral, so using $\mathbf{h}_{\text{eff},i}(u)$ instead of $\mathbf{h}_{\text{eff}}^* \{\bar{\mathbf{a}}_i(u), u\}$ yields the same influence function. This simplification is of great importance because otherwise, the optimal $\mathbf{h}_{\text{eff}}^*$ is only an interesting but impractical theoretical result. For most problems, computing $\mathbf{h}_{\text{eff}}^*$ is nearly impossible and would require extra modeling assumptions, which prevents the estimator from achieving the efficiency bound.

In our case, however, $\mathbf{h}_{\text{eff}}^*$ is simple to compute and the AIPW estimator achieves the optimal efficiency. A consistent estimate uses a sample version of (4) with IPW:

$$\widehat{\mathbf{h}}_{\text{eff},i}(u) = \frac{\sum_{j=1}^n I(\mathbf{q}_j = \mathbf{q}_i) \boldsymbol{\phi}(\mathbf{q}_j, x_j) Y_j(u) \delta_j / \widehat{G}(x_j)}{\sum_{j=1}^n I(\mathbf{q}_j = \mathbf{q}_i) Y_j(u) \delta_j / \widehat{G}(x_j)}. \quad (5)$$

Because $\mathbf{h}_{\text{eff},i}(u)$ is not a function of C_i , the independence between the censoring and the survival process gives

$$\begin{aligned} \mathcal{B}(\mathbf{h}_{\text{eff},i}, u) &= \frac{E\{\mathbf{h}_{\text{eff},i}(u) I(T_i \geq u) I(C_i \geq u)\}}{E\{I(T_i \geq u) I(C_i \geq u)\}} \\ &= \frac{E\{\mathbf{h}_{\text{eff},i}(u) Y_i(u)\}}{E\{Y_i(u)\}}. \end{aligned}$$

Therefore, we can approximate $\mathcal{B}(\mathbf{h}_{\text{eff}}, u)$ with

$$\widehat{\mathcal{B}}(\mathbf{h}_{\text{eff}}, u) = \frac{\sum_{i=1}^n \widehat{\mathbf{h}}_{\text{eff},i}(u) Y_i(u)}{Y(u)},$$

which satisfies

$$\sum_{i=1}^n \int \frac{\lambda^c(u) Y_i(u)}{\widehat{G}(u)} \{\widehat{\mathbf{h}}_{\text{eff},i}(u) - \widehat{\mathcal{B}}(\widehat{\mathbf{h}}_{\text{eff}}, u)\} du = 0.$$

This enables us to obtain the optimal AIPW estimator $\widehat{F}(t)$ by solving

$$\sum_{i=1}^n \left[\frac{\delta_i \boldsymbol{\phi}(\mathbf{q}_i, x_i)}{\widehat{G}(x_i)} + \int \frac{dN_i^c(u)}{\widehat{G}(u)} \{\widehat{\mathbf{h}}_{\text{eff},i}(u) - \widehat{\mathcal{B}}(\widehat{\mathbf{h}}_{\text{eff}}, u)\} \right] = 0. \quad (6)$$

The AIPW estimator is very easy to implement, especially compared with many other nonparametric or semiparametric problems, where the efficient estimator often involves additional model assumptions (Tsiatis and Ma 2004), solving integral equations (Rabinowitz 2000), and iterative procedures (Zhang, Tsiatis, and Davidian 2008).

In Section S.3 of the supplementary materials, we demonstrate that the AIPW estimator indeed has the efficient influence function (EFF), which corresponds to replacing $h(\cdot)$ with $\mathbf{h}_{\text{eff},i}(u)$ in (4). The variance of the efficient estimator is

$$V_{\text{eff}} = \text{cov}\{\boldsymbol{\phi}(\mathbf{Q}_i, T_i)\} + E \int \frac{\mathcal{B}\{(\boldsymbol{\phi} - \mathbf{h}_{\text{eff}})^{\otimes 2}, u\}}{G^2(u)} \lambda^c(u) Y_i(u) du,$$

which is estimated consistently by

$$\widehat{V}_{\text{eff}} = n^{-1} \sum_{i=1}^n \frac{\delta_i \boldsymbol{\phi}(\mathbf{q}_i, x_i) \boldsymbol{\phi}^T(\mathbf{q}_i, x_i)}{\widehat{G}(x_i)} + n^{-1} \sum_{i=1}^n \int \frac{\widehat{\mathcal{B}}_1\{(\boldsymbol{\phi} - \widehat{\mathbf{h}}_{\text{eff}})^{\otimes 2}, u\}}{\widehat{G}^2(u)} dN_i^c(u).$$

3.1.3 Subgroup-Specific Censoring. The IPW estimator (3) and the AIPW estimator (6) are designed for the case when the censoring distribution $G(\cdot)$ is common for all subjects in m subgroups. When this is not the case, subgroup-specific censoring distributions, $\widetilde{G}_j(t)$, $j = 1, \dots, m$, should be used. Specifically, $\widehat{G}(t)$ is replaced by the subgroup-specific Kaplan–Meier estimators

$$\widetilde{G}_j(t) = \prod_{\substack{x_i \leq t \\ \mathbf{q}_i = \mathbf{u}_j}} \left\{ 1 - \frac{1 - \delta_i}{\sum_{\mathbf{q}_k = \mathbf{u}_j} I(x_k \geq x_i)} \right\},$$

for $j = 1, \dots, m$. Consequently, the IPW estimating equation (3) changes to

$$n^{-1} \sum_{j=1}^m \sum_{i: \mathbf{q}_i = \mathbf{u}_j} \frac{\delta_i \boldsymbol{\phi}(\mathbf{q}_i, x_i)}{\widetilde{G}_j(x_i)} = 0,$$

and the corresponding estimated variance \widehat{V}_{ipw} changes analogously, with summation $\sum_{j=1}^m \sum_{i: \mathbf{q}_i = \mathbf{u}_j}$ replacing $\sum_{i=1}^n$, and a subgroup-specific

$$\widehat{B}_{1,j}(\mathbf{h}, u) = \frac{1}{\#\{i : \mathbf{q}_i = \mathbf{u}_j\} \widehat{S}_j(u)} \sum_{i: \mathbf{q}_i = \mathbf{u}_j} \frac{\delta_i \mathbf{h}(\mathbf{q}_i, x_i, \delta_i) I(x_i \geq u)}{\widetilde{G}_j(x_i)}$$

replacing the pooled $\widehat{B}_1(\mathbf{h}, u)$.

Similar changes apply to the AIPW estimator. Specifically, in Equations (5) and (6) and in the expression of $\widehat{\mathbf{V}}_{\text{eff}}$, we replace $\sum_{i=1}^n$ with $\sum_{j=1}^m \sum_{i:q_i=u_j}$, $\widehat{G}(t)$ with $\widetilde{G}_j(t)$, and $\widehat{\mathcal{B}}(\mathbf{h}_{\text{eff}}, u)$ with

$$\widehat{\mathcal{B}}_j(\mathbf{h}_{\text{eff}}, u) = \frac{\sum_{i:q_i=u_j} \mathbf{h}_{\text{eff},i}(u) Y_i(u)}{\sum_{i:q_i=u_j} Y_i(u)}.$$

It is worth noting that if we erroneously treat the censoring distribution as common when in fact it is not, the IPW estimator will not be consistent any more because the corresponding influence function no longer has mean zero. On the other hand, the AIPW estimator will still be consistent, although less efficient. This is a direct consequence of the double robustness property of the AIPW estimator, in that the validity of the complete-data influence function $\phi(\mathbf{q}, t)$ alone guarantees consistency of the AIPW estimator. However, since the efficiency of AIPW is achieved when the correct censoring model is used, treating the censoring distribution as identical across subgroups when they are not leads to efficiency loss. The issue of subgroup-specific censoring and the performance of IPW, AIPW, and their modified versions are illustrated in simulation studies in Section 4.3.

3.2 An IMP Estimator

Lipsitz, Ibrahim, and Zhao (1999) proposed a conditional estimating equation for regression with missing covariates by conditioning the complete-data estimating equation on the observed data. Similarly, with censored observations, we replace the unknown complete-data influence function with its conditional expectation, given that the event happens after the observed censoring time. Doing so yields the following imputed estimating equation:

$$\begin{aligned} 0 &= \sum_{i=1}^n [\delta_i \phi(\mathbf{q}_i, x_i) + (1 - \delta_i) E\{\phi(\mathbf{Q}_i, T_i) | T_i > x_i, \mathbf{q}_i\}] \\ &= \sum_{i=1}^n \{\delta_i \phi(\mathbf{q}_i, x_i) + (1 - \delta_i) \mathbf{h}_{\text{eff},i}(x_i)\}. \end{aligned}$$

In practice, with $\widehat{\mathbf{h}}_{\text{eff},i}(\cdot)$ as in (5), we obtain the IMP estimator by solving

$$0 = \sum_{i=1}^n \{\delta_i \phi(\mathbf{q}_i, x_i) + (1 - \delta_i) \widehat{\mathbf{h}}_{\text{eff},i}(x_i)\}.$$

While in many cases, the imputation method could lead to bias if the model of the missingness is misspecified, it is straightforward to see that our proposed IMP estimator is always consistent. In practice, we often, but not always, observe that it performs competitively in comparison with the optimal AIPW estimator. For inferences, we derive the influence function of the IMP estimator in Section S.4 of the supplementary materials and show that it has a complex form containing nested conditional expectations, and hence, is hardly useful practically. Asymptotic analysis for imputation-based estimation is often complex and can be rather involved even in parametric imputation procedures (Wang and Robins 1998; Robins and Wang 2000), which partially explains why the bootstrap method is usually favored in its inference.

4. SIMULATIONS

4.1 Simulation Design

We conducted comprehensive Monte Carlo simulations to illustrate the finite-sample performance of four groups of estimators, yielding a total of 11 different estimators. The first three groups of estimators include the IPW, optimal AIPW, and IMP estimators based on (a) the complete-data ordinary least-square (OLS) influence function, (b) the complete-data weighted least-square (WLS) influence function using the inverse of the variance as weights, and (c) the efficient (EFF) influence function; see the Appendix for the exact forms of the influence functions. The fourth group of estimators contains the two NPMLEs.

The primary goal of the simulation studies is to compare the bias and efficiency of the 11 estimators of the distribution function of an outcome subject to censoring in each genotype population without directly observing the genotypes. In the first two simulations, the distribution function $F(t)$ is a two-dimensional vector (i.e., $p = 2$ subpopulations). In the first simulation, we set $F_1(t) = \{1 - \exp(-t/4)\} / \{1 - \exp(-2.5)\}$ on the interval $(0, 10)$ and $F_2(t) = F_1(t)^{0.98}$ on the interval $(0, 5)$. In the second simulation, we set $F_1(t) = \{[1 - \exp(-t/4)] / [1 - \exp(-2.5)]\}^{0.5}$ on $(0, 10)$ and $F_2(t) = \{1 - \exp(-t/2)\} / \{1 - \exp(-2.5)\}$ on $(0, 5)$. Thus, the distributions in the first simulation belong to the proportional hazards model family, while they do not in the second. In both simulations, we let each random mixture proportion \mathbf{q}_i be one of $m = 4$ different possible vector values: $(1, 0)^T$, $(0.6, 0.4)^T$, $(0.2, 0.8)^T$, and $(0.16, 0.84)^T$. Our sample size was 500, and we generated a uniform censoring distribution to achieve moderate (20%) and high (50%) censoring rates.

Our third simulation mimics the COHORT data. With $p = 2$, we set $F_1(t) = [1 + \exp\{- (t - 63)/7\}]^{-0.9} / 0.995$ on $(0, 100)$, and $F_2(t) = 0.0007t$ on $(0, 53]$ and $F_2(t) = 0.022 + [1 + \exp\{- (t - 68)/7.5\}]^{-2}$ on $(53, 100)$. These distributions roughly mimic the estimated cumulative risk of death for HD gene mutation carriers and noncarriers, respectively, in the COHORT study (Figure 1). Analogous to the COHORT study, we used sample size $n = 4500$, generated $m = 6$ different mixture proportions: $(0, 1)^T$, $(0.5, 0.5)^T$, $(0.97, 0.03)^T$, $(0.75, 0.25)^T$, $(0.25, 0.75)^T$, and $(1, 0)^T$; and censored 65% of the observations with a uniformly distributed censoring process.

For each of the three simulations under different censoring rates, we evaluated all 11 estimators at different t -values. First, we ran 1000 Monte Carlo simulations to evaluate the pointwise bias, defined as $\widehat{F}(t) - F(t)$, at $t = 2.5$ in the first simulation (Table 1), at $t = 1.5$ in the second simulation (Table S.1, supplementary materials), and at $t = 70$ in the third simulation (Table S.2). The corresponding estimated standard errors for the IPW and AIPW estimators were based on $\widehat{\mathbf{V}}_{\text{ipw}}$ and $\widehat{\mathbf{V}}_{\text{eff}}$ given in Section 3.1.1 and Section 3.1.2, respectively, whereas bootstrap estimates were used to quantify the variability for the IMP estimator and the NPMLEs. Next, we evaluated the biases of the estimators across the entire range of t -values through an integrated absolute bias (IAB), defined as $\int_0^\infty |\widehat{F}_j(t) - F_j(t)| dt$, $j = 1, 2$, where $\widehat{F}_j(t)$ is the average estimated curves over multiple datasets. The results are in Table S.3 (upper half) and Tables S.2 and S.3 (upper half). In the first two simulations, the IAB was computed on a grid set

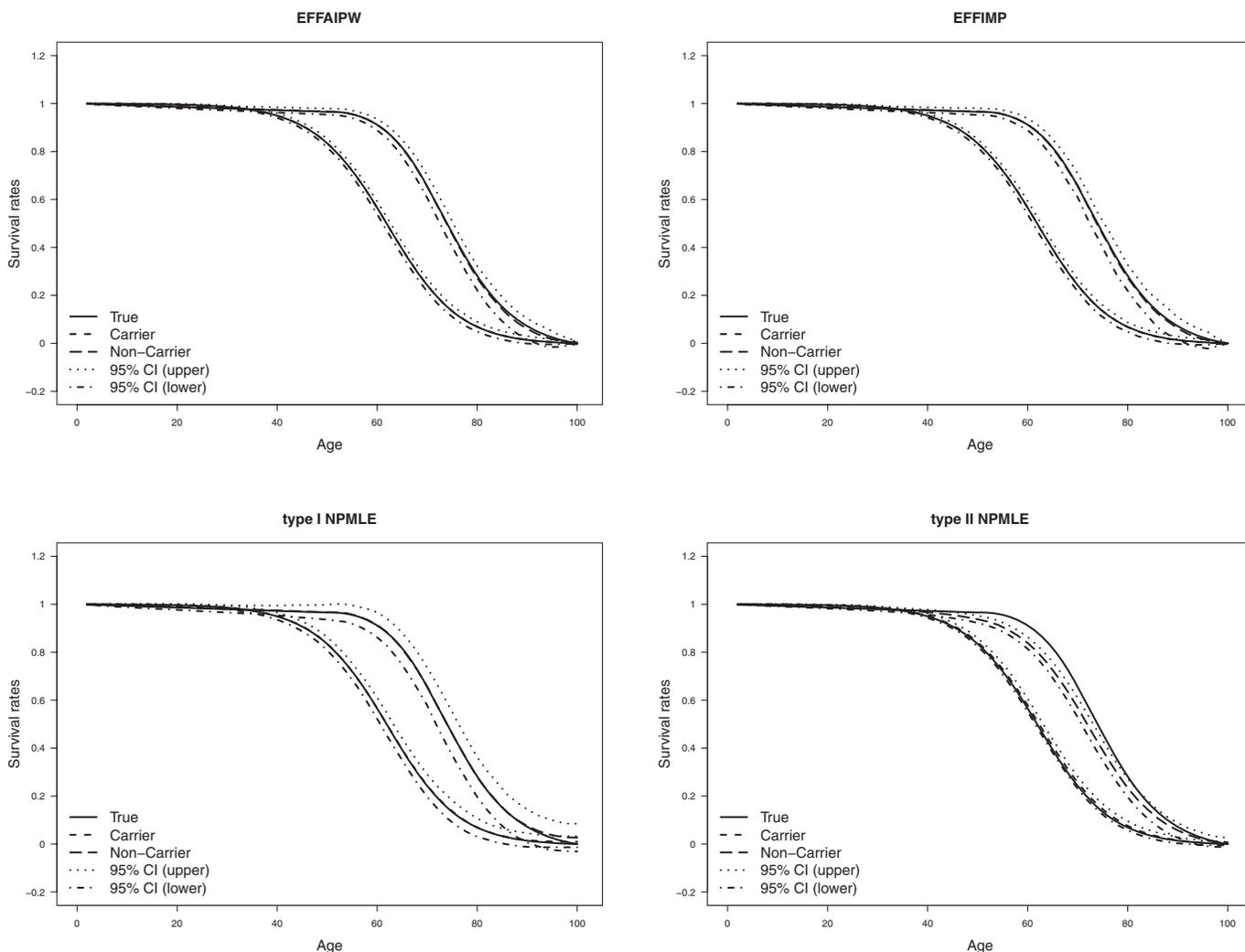


Figure 1. Simulation 3. True survival curve (solid) and the mean of 250 simulations at each time point (short-dashed for carrier group, long-dashed for noncarrier group), and 95% pointwise confidence band (upper band dotted, lower band dash-dotted) of the estimated survival curves. The mean and true survival curves are indistinguishable in the figures for the EFFIMP and EFFAIPW estimators. Sample size is 4500 and censoring rate is 65%.

with an increment of $\Delta = 0.1$ as $\sum_{i=1}^{100} |\bar{F}_1(x_i) - F_1(x_i)|\Delta$ and $\sum_{i=1}^{50} |\bar{F}_2(x_i) - F_2(x_i)|\Delta$, where $\bar{F}_j(x_i)$ ($j = 1, 2$) denotes the mean estimated distribution from 1000 datasets. In the third simulation, it was computed on a grid set with an increment of $\Delta = 2$ on $(0,100)$ as $\sum_{i=1}^{50} |\bar{F}_j(x_i) - F_j(x_i)|\Delta$ ($j = 1, 2$), where $\bar{F}_j(x_i)$ denotes the mean estimated distribution from 250 datasets.

4.2 Simulation Results

The results in Table 1 and Tables S.1 and S.2 indicate that all the nonparametric estimators we propose have ignorable finite-sample biases, while the Type II NPMLE has much larger bias. At high censoring rates, the bias for the Type II NPMLE is much greater and the coverage probability is much lower than the pre-specified nominal level. Moreover, the bias is not a finite-sample effect since even at a sample size of $n = 4500$, the bias persists. Despite its asymptotic consistency, the Type I NPMLE also shows substantial bias when the censoring rate increases. This is because in the estimation procedure of the Type I NPMLE,

the mixture nature of the model is not taken advantage of at the maximization step. The Kaplan–Meier estimation in some subgroups could be based on very small sample sizes, which can make the overall estimation unreliable.

Compared with the proposed estimators, the Type I NPMLE has, for the most part, larger estimation variability, and the increased variability is rather substantial for the $F_2(t)$ estimation. In particular, the WLS-based estimators have much less variability than the Type I NPMLE. Ma and Wang (2012) showed that the Type I NPMLE for uncensored data belongs to the WLS family with weights $w_i = 1/n_i$, where n_i is the number of subjects who share the same q_i . In other words, the Type I NPMLE essentially downweights individuals belonging to large subgroups. Since such a weighting strategy is highly undesirable, it is not surprising to see that the weights in the WLS provide an improvement over those in the Type I NPMLE.

In contrast, the three proposed nonparametric estimators have satisfactorily small biases and are more efficient compared with the Type I NPMLE. The optimal AIPW and IMP estimators both improve upon IPW in terms of estimation efficiency. When the

Table 1. Simulation 1. Bias, empirical standard deviation (Emp. SD), average estimated standard deviation (Est. SD), and 95% coverage (95% cov.) of 11 estimators

Estimator	$F_1(t) = 0.5063$				$F_2(t) = 0.5132$			
	Bias	Emp. SD	Est. SD	95% cov.	Bias	Emp. SD	Est. SD	95% cov.
Group 1: OLS-based, censoring rate = 20%								
IPW	0.0013	0.0424	0.0427	0.9460	-0.0026	0.0408	0.0418	0.9590
AIPW	0.0003	0.0396	0.0402	0.9450	-0.0021	0.0390	0.0393	0.9520
IMP	0.0004	0.0396	0.0401	0.9460	-0.0022	0.0389	0.0392	0.9500
Group 2: WLS-based, censoring rate = 20%								
IPW	0.0013	0.0424	0.0427	0.9450	-0.0026	0.0408	0.0418	0.9590
AIPW	0.0003	0.0396	0.0402	0.9440	-0.0021	0.0390	0.0393	0.9520
IMP	0.0004	0.0396	0.0401	0.9460	-0.0022	0.0389	0.0392	0.9500
Group 3: EFF-based, censoring rate = 20%								
IPW	0.0011	0.0427	0.0431	0.9520	-0.0029	0.0418	0.0433	0.9630
AIPW	0.0004	0.0399	0.0404	0.9480	-0.0022	0.0393	0.0399	0.9540
IMP	0.0004	0.0398	0.0403	0.9430	-0.0022	0.0391	0.0394	0.9500
Group 4: NPMLE, censoring rate = 20%								
Type I	0.0000	0.0462	0.0466	0.9470	0.0007	0.0881	0.0897	0.9230
Type II	-0.0135	0.0364	0.0364	0.9300	0.0075	0.0308	0.0313	0.9390
Group 1: OLS-based, censoring rate = 50%								
IPW	0.0043	0.0733	0.0683	0.9290	0.0004	0.0722	0.0668	0.9290
AIPW	0.0010	0.0452	0.0459	0.9530	-0.0026	0.0452	0.0448	0.9450
IMP	0.0043	0.0486	0.0496	0.9580	0.0012	0.0484	0.0480	0.9480
Group 2: WLS-based, censoring rate = 50%								
IPW	0.0044	0.0732	0.0683	0.9290	0.0004	0.0722	0.0668	0.9280
AIPW	0.0010	0.0452	0.0459	0.9520	-0.0026	0.0451	0.0448	0.9460
IMP	0.0043	0.0486	0.0496	0.9590	0.0012	0.0484	0.0480	0.9470
Group 3: EFF-based, censoring rate = 50%								
IPW	-0.0021	0.0740	0.0702	0.9290	0.0005	0.0799	0.0744	0.9300
AIPW	-0.0006	0.0465	0.0464	0.9500	-0.0008	0.0470	0.0455	0.9480
IMP	0.0031	0.0493	0.0508	0.9610	0.0026	0.0495	0.0489	0.9480
Group 4: NPMLE, censoring rate = 50%								
Type I	0.0011	0.0531	0.0537	0.9410	0.0001	0.1058	0.1030	0.9190
Type II	-0.0405	0.0407	0.0417	0.8350	0.0312	0.0381	0.0382	0.8750

NOTE: Sample size $n = 500$, 20% and 50% censoring rate, and 1000 simulations.

censoring rate is moderate, IMP and AIPW perform similarly, while when the censoring rate increases, the superiority of the optimal AIPW over IMP becomes more notable. The similarity of the results in the first three groups of estimators suggests that the estimation efficiency is not sensitive to the choice of the noncensored-data influence function ϕ . The same insensitivity of estimation efficiency to the choice of influence function is also observed for the noncensored-data case (Ma and Wang 2012). This phenomenon proves beneficial, especially in the censored-data analysis since Robins and Rotnitzky (1992) remarked that the best complete-data influence function does not necessarily yield an optimal censored-data influence function, and that finding the optimal member usually requires a computationally intensive procedure. Finally, the estimated standard error matches reasonably well with the empirical standard error, while the 95% coverage probability is close to the nominal level.

In Figure 1, we examine the entire estimated survival curve $1 - \hat{F}(t)$ in the third simulation, which mimics the COHORT study. In particular, we display results from the imputation estimator (EFFIMP) and the AIPW estimator (EFFAIPW), both

based on the complete-data EFF, as representatives of the proposed estimators, and compare them with the two NPMLEs. To evaluate these estimators, we plotted the true survival curves and the mean estimated survival curves from 250 simulated datasets along with the 95% pointwise confidence bands. The EFFIMP and EFFAIPW estimators perform satisfactorily throughout the entire range of t , while the Type I NPMLE starts to exhibit large variability and small-sample estimation bias as time progresses. This confirms our previous observation that the Type I NPMLE suffers from the small-subgroup sample size difficulty and the instability of the Kaplan–Meier estimation procedure near the maximum event time. The Type II NPMLE also shows a nonignorable bias for a wide range of t 's.

To illustrate the overall bias of the 11 estimators across the entire range of t -values, we further provide the IAB for all three simulations in Table 2 (upper half) and Tables S.2 and S.3 (upper half). Nearly all estimators have very small IAB, whereas the Type II NPMLE can yield a bias 10 times larger than the other estimators. For the same estimator in each simulation, the IAB also tends to increase with higher censoring rate.

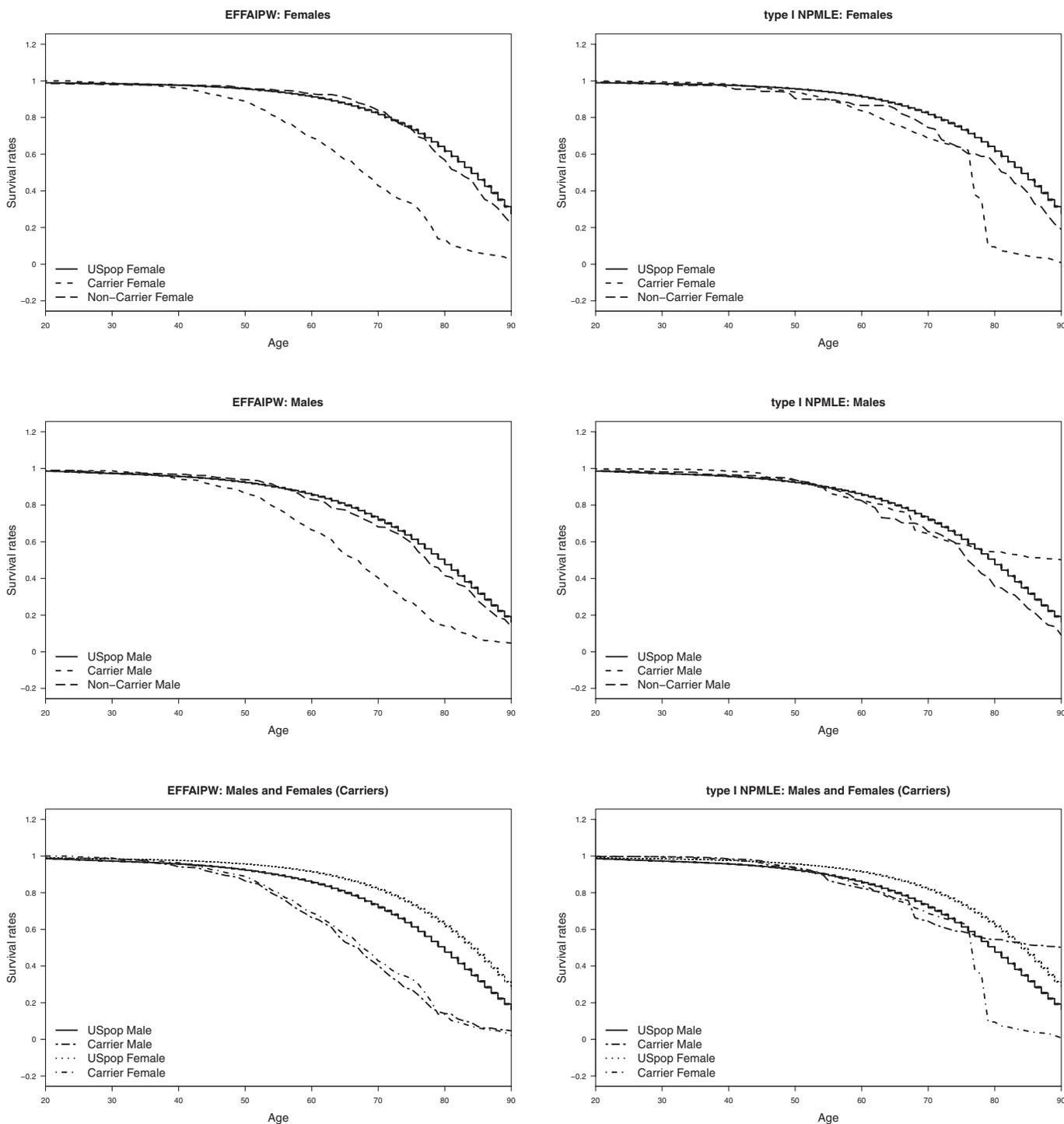


Figure 2. Estimated survival curves for the COHORT data stratified by gender using the EFFAIPW estimator (left) and the Type I NPMLE (right). The curves “USpop Male” and “USpop Female” correspond to Kaplan–Meier estimated survival rates for the general male and female U.S. populations in 2003, respectively. Bottom two figures compare the estimated survival curves for the male and female HD gene mutation carriers with that of the general male and female U.S. populations using the EFFAIPW estimator (left) and the Type I NPMLE (right).

4.3 Subgroup-Specific Censoring

We now examine the performance of the original IPW and AIPW estimators proposed in Sections 3.1.1 and 3.1.2, respectively, as well as the performance of their modified versions in Section 3.1.3, when the true censoring distribution is different across different subgroups. We extend the first simulation by generating the censoring times from the proportional hazards

distribution

$$G(t | \mathbf{q}_i) = 1 - \exp\{-\gamma_1 t^{\gamma_2} \exp(\gamma_3)\},$$

where we set $\gamma_3 = 0$ if $\mathbf{q}_i = (1, 0)^T$ or $\mathbf{q}_i = (0.6, 0.4)^T$, and set $\gamma_3 = 0.2$ if $\mathbf{q}_i = (0.2, 0.8)^T$ or $\mathbf{q}_i = (0.16, 0.84)^T$. Here, γ_1, γ_2 remain the same across the different subgroups and were chosen to achieve, respectively, moderate (20%) and high (50%)

Table 2. Simulation 1. Integrated absolute bias (IAB)

Estimator	Censoring rate			
	20%		50%	
	$F_1(t)$	$F_2(t)$	$F_1(t)$	$F_2(t)$
Group 1: OLS-based				
IPW	0.0103	0.0060	0.0332	0.0096
AIPW	0.0091	0.0055	0.0305	0.0077
IMP	0.092	0.0056	0.0348	0.0146
Group 2: WLS-based				
IPW	0.0102	0.0060	0.0345	0.0095
AIPW	0.0370	0.0055	0.1062	0.0092
IMP	0.0092	0.0056	0.0357	0.0145
Group 3: EFF-based				
IPW	0.0095	0.0062	0.0170	0.0106
AIPW	0.0090	0.0053	0.0462	0.0070
IMP	0.0090	0.0055	0.0325	0.0195
Group 4: NPMLE				
Type I	0.0104	0.0116	0.0174	0.0337
Type II	0.0996	0.0374	0.2483	0.1367
Group 1: OLS-based				
IPW [†]	0.0922	0.0632	0.0651	0.0412
AIPW [†]	0.0110	0.0045	0.0317	0.0101
IPW*	0.0159	0.0052	0.0181	0.0073
AIPW*	0.0148	0.0052	0.0240	0.0079
Group 2: WLS-based				
IPW [†]	0.0923	0.0632	0.0658	0.0413
AIPW [†]	0.0517	0.0045	0.1024	0.0111
IPW*	0.0169	0.0051	0.0218	0.0074
AIPW*	0.0591	0.0051	0.1031	0.0092
Group 3: EFF-based				
IPW [†]	0.1104	0.0791	0.0642	0.0495
AIPW [†]	0.0122	0.0054	0.0392	0.0094
IPW*	0.0151	0.0053	0.0149	0.0071
AIPW*	0.0112	0.0073	0.0351	0.0068

NOTES: Upper half of the table: the true censoring distribution is independent of \mathbf{Q}_i ; $G(t)$ is estimated using a common Kaplan–Meier estimator of the censoring distribution. Lower half of the table: the true censoring distribution is subgroup-specific; $G(t)$ is estimated using a common Kaplan–Meier estimator (denoted by [†]) or a subgroup-specific Kaplan–Meier estimator (denoted by *).

censoring proportions. The survival times and \mathbf{q}_i were generated as in the first simulation in Section 4.1. We first implemented the original IPW and AIPW estimators, which ignore the subgroup-specific censoring pattern and simply use a pooled censoring distribution $\widehat{G}(t)$. We then implemented the modified IPW and AIPW estimators, which incorporate the subgroup-specific censoring distributions by obtaining $\widetilde{G}_j(t)$ as described in Section 3.1.3. As in our earlier analyses, we investigated the pointwise bias at $t = 2.5$, as well as the IAB across the entire range of t .

Table 3 shows the pointwise bias of the IPW and AIPW estimators when using a pooled estimated censoring distribution and a subgroup-specific censoring distribution. It is evident that the IPW has substantial bias if the pooled censoring estimate is used, indicating that the original IPW is not applicable when the censoring is subgroup-specific. However, its bias is substantially reduced as soon as the subgroup-specific censoring is taken into account. The AIPW, on the other hand, is robust to misspeci-

fication and has small bias regardless of whether a pooled or a subgroup-specific censoring estimate is used. For the AIPW, ignoring the subgroup-specific censoring pattern only incurs a small efficiency loss, mostly for the OLS-based and WLS-based estimators. For the EFFAIPW, the efficiency loss is minimal.

The IAB (lower half of Table 2) further indicates that when the pooled censoring distribution is used, the IPW has much larger overall bias than the AIPW. However, after incorporating the subgroup-specific censoring distribution, the modified IPW and AIPW have similar magnitudes of the IAB. In a separate analysis where we extended the second simulation in Section 4.1, we also found similar behaviors for the pointwise bias (Table S.4) and IAB (lower half of Table S.3).

5. ANALYSIS OF COHORT DATA

Data from the COHORT study consist of 4587 relatives of the proband participants who have different mixing proportions for being carriers or noncarriers of the HD gene mutation. Computation of these mixing proportions is discussed in section 2 of Wacholder et al. (1998) and in Wang et al. (2008). The event time of interest is age of death, and roughly 68% of the data is censored. A main research interest is estimating the age-at-death distribution or the survival function for carriers and noncarriers to assess the effects of HD mutation on survival. The severity of HD warrants that noncarriers tend to live longer, so we expect to see lower survival rates for the carrier group.

Since it is well known that survival rates differ by gender, we stratified the COHORT data by gender (2367 males and 2220 females) and analyzed the effects of HD gene mutation on the male and female subpopulations. The quantity of interest, $\mathbf{1} - \mathbf{F}(t)$, is a four-dimensional vector (i.e., $p = 4$), where $1 - F_1(t)$ and $1 - F_2(t)$ denote the survival functions for male noncarriers and carriers, respectively, and $1 - F_3(t)$ and $1 - F_4(t)$ denote analogous functions for females. Furthermore, the mixture proportions \mathbf{Q}_i are four-dimensional vectors with the first two components corresponding to mixture proportions for male noncarriers and carriers, and the last two components for female noncarriers and carriers, respectively. To estimate $\mathbf{1} - \mathbf{F}(t)$, we implemented several theoretically consistent nonparametric estimators, including the Type I NPMLE and the complete-data EFFAIPW as representatives of the already existing and newly proposed methods, respectively.

To examine the performance of these estimators, we first compare the results for the male and female noncarrier groups with the general male and female U.S. populations in 2003 (Arias 2006). These survival rates should be similar since the risk in noncarriers for both genders would reflect the general population if there is minimal ascertainment bias in family members. Figure 2 and Tables 4 and 5 (lowest panel) indicate that the EFFAIPW outperforms the Type I NPMLE in capturing the behavior of the general male and female U.S. populations. In fact, comparing the noncarrier female estimates and general female population, the EFFAIPW has an IAB less than half of that of the Type I NPMLE. Likewise, for the noncarrier males, the EFFAIPW has an IAB about half of that of the Type I NPMLE. Hence, the EFFAIPW appears to be a more reasonable estimator for analysis.

Table 3. Simulation 1. Bias, empirical standard deviation (Emp. SD), average estimated standard deviation (Est. SD), and 95% coverage (95% cov.) of 11 estimators

Estimator	$F_1(t) = 0.5063$				$F_2(t) = 0.5132$			
	Bias	Emp. SD	Est. SD	95% cov.	Bias	Emp. SD	Est. SD	95% cov.
Group 1: OLS-based, censoring rate = 20%								
IPW [†]	-0.0132	0.0439	0.0443	0.9440	0.0131	0.0433	0.0428	0.9360
AIPW [†]	0.0007	0.0386	0.0399	0.9520	-0.0018	0.0380	0.0381	0.9490
IPW*	0.0019	0.0395	0.0395	0.9370	-0.0012	0.0384	0.0389	0.9530
AIPW*	0.0018	0.0391	0.0391	0.9430	-0.0012	0.0383	0.0384	0.9490
Group 2: WLS-based, censoring rate = 20%								
IPW [†]	-0.0132	0.0439	0.0443	0.9440	0.0131	0.0433	0.0428	0.9360
AIPW [†]	0.0007	0.0386	0.0399	0.9520	-0.0018	0.0380	0.0381	0.9490
IPW*	0.0019	0.0395	0.0395	0.9400	-0.0012	0.0384	0.0389	0.9530
AIPW*	0.0018	0.0390	0.0391	0.9430	-0.0012	0.0383	0.0384	0.9500
Group 3: EFF-based, censoring rate = 20%								
IPW [†]	-0.0153	0.0459	0.0452	0.9350	0.0163	0.0446	0.0450	0.9360
AIPW [†]	0.0004	0.0392	0.0401	0.9470	-0.0014	0.0387	0.0386	0.9470
IPW*	0.0020	0.0392	0.0396	0.9390	-0.0014	0.0384	0.0393	0.9550
AIPW*	0.0013	0.0393	0.0393	0.9390	-0.0007	0.0387	0.0390	0.9500
Group 1: OLS-based, censoring rate = 50%								
IPW [†]	-0.0077	0.0708	0.0682	0.9320	0.0130	0.0729	0.0663	0.9230
AIPW [†]	0.0014	0.0448	0.0472	0.9570	-0.0027	0.0463	0.0442	0.9380
IPW*	0.0032	0.0502	0.0474	0.9410	-0.0028	0.0492	0.0477	0.9380
AIPW*	0.0021	0.0473	0.0449	0.9410	-0.0021	0.0471	0.0455	0.9330
Group 2: WLS-based, censoring rate = 50%								
IPW [†]	-0.0076	0.0707	0.0682	0.9310	0.0130	0.0729	0.0663	0.9220
AIPW [†]	0.0014	0.0448	0.0472	0.9550	-0.0027	0.0462	0.0442	0.9380
IPW*	0.0033	0.0504	0.0474	0.9390	-0.0028	0.0492	0.0477	0.9370
AIPW*	0.0022	0.0473	0.0449	0.9420	-0.0021	0.0471	0.0455	0.9330
Group 3: EFF-based, censoring rate = 50%								
IPW [†]	-0.0148	0.0753	0.0707	0.9250	0.0169	0.0794	0.0736	0.9240
AIPW [†]	0.0008	0.0463	0.0477	0.9570	-0.0021	0.0487	0.0449	0.9350
IPW*	0.0020	0.0481	0.0466	0.9480	-0.0024	0.0485	0.0477	0.9400
AIPW*	0.0003	0.0472	0.0454	0.9470	-0.0003	0.0483	0.0462	0.9390s

NOTES: The censoring distribution is subgroup-specific, and $G(t)$ is estimated using a common Kaplan–Meier estimator (denoted by [†]) or a subgroup-specific Kaplan–Meier estimator (denoted by *). Sample size $n = 500$, 20% and 50% censoring rate, and 1000 simulations.

The EFFAIPW depicted in Figure 2 shows a steep difference in the estimated survival rates between the carrier and the noncarrier groups for both genders. In addition, the bottom-left figure suggests that male carriers tend to have only slightly lower survival rates than female carriers. For example, at age 65, male carriers have a cumulative risk of death of 46.9% (95% CI: 40.9%–53%) whereas female carriers have a cumulative risk of death of 43.0% (95% CI: 37.4%–48.7%). This slight difference, in combination with the overlapping 95% confidence band (not shown here), suggests that HD mutation affects males and females equally. The observed lack of gender effects from EFFAIPW agrees with some earlier studies that also did not find a gender effect in either the mean survival times of HD patients (Harper 1996) or the progression of HD (Marder et al. 2000). In contrast, the Type I NPMLE suggests that male carriers have much better survival rates than female carriers, and sometimes, even slightly better survival rates than noncarriers—a behavior contradictory to the existing clinical literature.

The upper panel of Table 5 further presents the area under the survival curves, which can be interpreted as the expected years of

life. Hence, the difference in the area under two survival curves represents the expected years of life lost for one compared with the other. Based on the EFFAIPW, the estimated expected years of life lost for mutation carriers compared with noncarriers is 9.06 years in males and 12.76 in females. In contrast, the Type I NPMLE estimates longer expected years of life for male carriers, which is unreasonable.

Our further investigation reveals that the poor performance of the Type I NPMLE on the COHORT data is partially due to small-sample sizes in several Q_i subgroups. When we remove 211 subjects pertaining to three subgroups with sample sizes no more than 3% of the data, the behavior of the Type I NPMLE improves (Figure S.1 in the supplementary materials). The improvement is reflected in both the IAB between noncarriers and the U.S. population and the survival rates for carriers in both genders. However, there is still a large difference in the noncarrier female Type I NPMLE estimate in terms of expected life lost compared with the U.S. population (Table 5, middle panel).

Lastly, in Table 6, we present the estimated conditional probabilities of surviving the next several years in five-year intervals

Table 4. Estimated survival rates and 95% confidence intervals (in parentheses) based on EFFAIPW and Type I NPMLE for carrier and noncarrier groups in the COHORT data stratified by gender

Age	USpop	Noncarrier		Carrier	
		EFFAIPW	Type I NPMLE	EFFAIPW	Type I NPMLE
Males					
30	97.2 (97.1–97.3)	97.5 (96.6–98.3)	98.1 (96.8–99.4)	98.7 (98.2–99.3)	99.7 (99.4–99.9)
35	96.5 (96.4–96.6)	97.4 (96.6–98.3)	97.9 (96.8–99.0)	97.2 (96.2–98.2)	99.4 (98.9–99.9)
40	95.6 (95.5–95.7)	96.8 (95.7–98.0)	96.5 (94.2–98.7)	94.1 (92.4–95.8)	98.4 (97.6–99.2)
45	94.3 (94.1–94.4)	95.5 (94.0–97.1)	95.1 (92.1–98.1)	91.2 (89.0–93.4)	97.4 (96.1–98.7)
50	92.3 (92.1–92.4)	93.9 (92.0–95.8)	93.8 (90.0–97.5)	86.4 (83.6–89.3)	93.4 (88.8–98.1)
55	89.4 (89.2–89.6)	89.7 (86.9–92.5)	88.9 (83.5–94.3)	77.9 (74.0–81.9)	86.5 (78.8–94.3)
60	85.5 (85.3–85.7)	83.2 (79.5–87.0)	82.4 (76.1–88.7)	66.5 (61.4–71.5)	82.4 (72.4–92.5)
65	79.9 (79.6–80.1)	77.4 (72.9–81.9)	72.6 (65.8–79.3)	53.1 (47.0–59.1)	76.9 (65.8–88.0)
70	72.0 (71.7–72.3)	68.1 (62.7–73.5)	65.4 (56.3–74.5)	40.3 (33.6–47.0)	64.4 (43.5–85.3)
75	61.3 (61.0–61.6)	59.4 (53.3–65.6)	54.8 (45.5–64.2)	27.3 (20.5–34.0)	58.9 (39.3–78.5)
Females					
30	98.4 (98.4–98.5)	98.0 (97.1–98.9)	98.7 (98.1–99.3)	98.9 (98.1–99.7)	99.5 (99.1–99.9)
35	98.1 (98.0–98.2)	97.8 (96.8–98.7)	97.5 (95.6–99.4)	97.7 (96.6–98.9)	99.0 (98.4–99.5)
40	97.6 (97.5–97.7)	97.7 (96.7–98.8)	96.5 (93.6–99.4)	96.3 (94.8–97.8)	98.2 (97.5–98.9)
45	96.7 (96.6–96.9)	97.4 (96.2–98.6)	94.3 (89.8–98.7)	92.6 (90.5–94.7)	96.2 (95.3–97.1)
50	95.5 (95.4–95.7)	96.1 (94.4–97.8)	90.3 (85.0–95.7)	88.9 (86.2–91.7)	93.8 (92.4–95.3)
55	93.8 (93.6–94.0)	94.7 (92.7–96.7)	89.5 (84.1–94.9)	79.9 (76.3–83.5)	89.9 (88.2–91.6)
60	91.2 (91.1–91.4)	93.0 (90.3–95.6)	86.6 (80.5–92.7)	69.1 (64.5–73.8)	83.7 (81.6–85.8)
65	87.3 (87.1–87.5)	91.0 (87.6–94.4)	84.9 (79.2–90.7)	57.0 (51.3–62.6)	75.9 (73.3–78.5)
70	81.6 (81.3–81.8)	84.0 (79.3–88.6)	74.5 (67.3–81.6)	42.8 (36.5–49.2)	68.7 (66.2–71.1)
75	73.3 (73.0–73.6)	73.9 (68.0–79.8)	62.8 (55.0–70.6)	33.2 (26.6–39.9)	63.6 (60.7–66.5)

NOTE: Survival rates are compared with Kaplan–Meier estimated survival rates for the general male and female U.S. populations (USpop) in 2003.

for a subject alive at a given age. These probabilities are based on the EFFAIPW and are stratified by gender and mutation status. For example, a 35-year-old female carrier has a 94.71% chance of surviving the next 10 years, and an 81.76% chance of surviving

the next 20 years, compared with a 99.63% and a 96.81% chance in the case of a female noncarrier, respectively. Such probabilities are useful for patients when interpreting mutation testing results, and may help them make lifetime decisions, such as having children.

In conclusion, using the more reliable EFFAIPW estimator, our analysis suggests that mutation carriers tend to have much lower survival rates than noncarriers, and the mutation equally affects survival rates in both genders. These survival rates are the first in the literature obtained from a sample of family members, and they highlight the deleterious effects of HD mutation on survival. The estimated survival rates in noncarriers closely resemble that of the U.S. population, which illustrates minimal ascertainment bias and reflects the advantage of analyzing family members whose information is not used in the initial recruitment of probands.

6. DISCUSSION

We propose two IPW-based estimators and an IMP estimator for censored mixture data, among which the AIPW achieves the optimal efficiency based on a given complete-data influence function. These estimators are easy to compute and do not involve any iterative procedures. When the sample size is small and the censoring rate is moderate, the IMP estimator can sometimes compete or even outperform the asymptotically optimal AIPW estimator. We also point out the surprising results of the inefficiency of the Type I NPMLE and the inconsistency of the Type II NPMLE proposed in the literature. Our finite-sample simulations suggest that the efficiency loss of the Type I NPMLE and the bias of the Type II estimator can be quite

Table 5. Summary statistics for the COHORT data and the general male and female U.S. populations (USpop) in 2003 from age 20 to 90 years

	Males		Females	
	Noncarrier	Carrier	Noncarrier	Carrier
Expected years of life (area under the survival curves)				
USpop	55.5559		59.8726	
EFFAIPW	54.8759	45.8229	59.6289	46.8695
Type I NPMLE	53.6082	57.7125	57.0451	52.1563
Type I NPMLE ^a	56.5764	47.3792	61.1127	47.2930
Expected years of life lost compared with the U.S. population				
EFFAIPW	0.6800	9.7330	0.2436	13.0031
Type I NPMLE	1.9477	-2.1566	2.8275	7.7162
Type I NPMLE ^a	-1.0205	8.1767	-1.2401	12.5796
IAB* between an estimator and the U.S. population				
EFFAIPW	1.4026	10.0315	1.2957	13.1922
Type I NPMLE	2.6079	3.9104	2.9016	8.1144
Type I NPMLE ^a	1.1071	8.7317	1.8435	12.7560

^aComputed under a subsample by removing subjects in Q_i groups with small sample sizes.
*Integrated Absolute Bias.

Table 6. Estimated conditional probabilities of survival for the COHORT data based on the EFFAIPW estimator

Current age	Conditional probability individual alive in the next span of years							
	5 years	10 years	15 years	20 years	5 years	10 years	15 years	20 years
	Male carriers				Male noncarriers			
30	0.9843	0.9533	0.9239	0.8753	0.9996	0.9937	0.9803	0.9636
35	0.9686	0.9387	0.8893	0.8017	0.9941	0.9808	0.9640	0.9206
40	0.9692	0.9182	0.8277	0.7061	0.9866	0.9698	0.9261	0.8594
45	0.9474	0.8540	0.7285	0.5818	0.9830	0.9387	0.8711	0.8099
50	0.9015	0.7690	0.6142	0.4666	0.9550	0.8862	0.8239	0.7251
55	0.8530	0.6813	0.5176	0.3498	0.9280	0.8628	0.7593	0.6628
60	0.7987	0.6068	0.4101	0.2122	0.9297	0.8182	0.7142	0.4984
	Female carriers				Female noncarriers			
30	0.9887	0.9744	0.9363	0.8997	0.9977	0.9975	0.9940	0.9806
35	0.9855	0.9471	0.9100	0.8176	0.9998	0.9963	0.9829	0.9689
40	0.9610	0.9233	0.8296	0.7177	0.9965	0.9831	0.9691	0.9511
45	0.9608	0.8633	0.7469	0.6153	0.9865	0.9725	0.9544	0.9346
50	0.8985	0.7773	0.6404	0.4814	0.9858	0.9674	0.9473	0.8737
55	0.8651	0.7127	0.5358	0.4156	0.9814	0.9610	0.8863	0.7802
60	0.8239	0.6194	0.4804	0.1971	0.9792	0.9031	0.7950	0.6116

substantial, and the finite-sample bias of the Type I NPMLE can be nonignorable when the subgroup sample size is small or the estimation region is close to the upper end of the distribution support. Caution should be applied in interpreting inconsistency of the Type II NPMLE, which occurs when a pure nonparametric model is used. Parametric models and semiparametric models, such as the Cox proportional hazards model with a nonparametric baseline or piecewise exponential model, are expected to be consistent (Zeng and Lin 2007).

In a special case when the data arise from a single distribution (i.e., $p = 1$), the IPW, AIPW, IMP, and the two NPMLEs are all equivalent to the familiar Kaplan–Meier estimator. This indicates the complexity arising from the mixture nature. Through extensive simulation studies, we demonstrate that the proposed AIPW has satisfactory small bias and is more efficient than the Type I NPMLE even when the censoring rate is high.

The optimal AIPW estimator also provides reasonable survival rate estimates for both genders and different mutation status in the COHORT study. Since genetic testing of HD mutation is commercially available, the estimated survival rates are useful in genetic counseling when a subject, with a family history of HD, needs to decide on whether to undergo genetic testing. Understanding the mortality rates associated with a positive testing result may make the subject more inclined to determine his/her mutation status and seek treatments. In addition, in a future work, it may be of interest to estimate the survival rates as a function of the number of CAG repeats in carriers.

In some genetic studies, the relatives are included through their probands, and there might be potential ascertainment bias. If the HD gene mutation carriers' probands are randomly sampled from the population of all carriers, then the estimation from the relatives can be generalized to the population of all carriers. This corresponds to no ascertainment bias. However, when there is heterogeneity in the survival function of HD gene mutation carriers (e.g., there exists another gene influencing the survival function) and if the carriers' probands are not a representative sample of the population of all HD mutation carriers, then estimation based on their relatives may be biased (Begg 2002).

For example, if there is a second gene that decreases survival in HD mutation carriers, then oversampling of probands with the second gene may lead to an upward bias of the risk of death, and undersampling would lead to a downward bias. However, in the analysis of the COHORT data, the estimated survival function in noncarriers is reasonably close to the general population estimates obtained from the census data. This is an indication that the COHORT sample is not likely to be subject to severe ascertainment bias. Otherwise, the noncarrier distribution estimated from the COHORT relative data would be very different from the that for the general population due to the distorted distribution of additional risk factors.

Finally, since the survival distribution is well known to be different between two genders in the general population, we carried out the COHORT analysis separately for each gender. It may be desirable to test the gender difference among the HD gene mutation carriers/noncarriers. This amounts to testing $H_0 : \mathbf{F}^1(t) = \mathbf{F}^2(t)$ at all t versus $H_1 : \mathbf{F}^1(t) \neq \mathbf{F}^2(t)$ for at least one t , where $\mathbf{F}^1(\cdot)$ is the vector of cumulative distribution functions for males and $\mathbf{F}^2(\cdot)$ for females. Among various methods of performing the test, a convenient choice is permutation. Specifically, we compute the test statistic $v_0 = \sup_t \|\hat{\mathbf{F}}^1(t) - \hat{\mathbf{F}}^2(t)\|$ from the observed data, where $\|\cdot\|$ is the L_2 -norm. Since under the null hypothesis, the two genders have identical distributions, we can randomly permute the gender variable to obtain a permuted sample. Perform such a permutation B times for some large B , and recompute the test statistic v_b based on the b th permuted sample, $b = 1, \dots, B$. The p -value is then $\sum_{b=1}^B I(v_b \geq v_0)/B$. If interest only lies in the gender difference in the carrier population, one may extract the corresponding component in $\mathbf{F}^1(t)$ and $\mathbf{F}^2(t)$ to perform the test.

APPENDIX: INFLUENCE FUNCTION OF CONSISTENT ESTIMATORS WITH COMPLETE DATA

When there is no censoring [i.e., $\delta_i = 1$ for all i in (1)], Ma and Wang (2012) adopted a pure nonparametric model of the genotype-specific distributions without assuming any parametric form of the

density function. They proposed a general class of consistent estimators and identified the efficient member of the class. The complete set of all influence functions of the consistent estimators for $F(t)$ can be characterized as (Ma and Wang 2012)

$$\left\{ \phi(\mathbf{q}, s) : \phi(\mathbf{q}, s) = \mathbf{d}(\mathbf{q}, s) - F(t) - \mathbf{B}\mathbf{1}_p, \int \mathbf{d}(\mathbf{Q}, s) \mathbf{Q}^T p_{\mathbf{Q}}(\mathbf{Q}) d\mu(\mathbf{Q}) = I(s \leq t) \mathbf{I}_p + \mathbf{B} \right\},$$

where \mathbf{I}_p is a p -dimensional identity matrix, $\mathbf{d}(\mathbf{q}, s)$ is a vector of real functions (for qualified choices of $\mathbf{d}(\mathbf{q}, s)$, see Ma and Wang 2012), \mathbf{B} is an arbitrary $p \times p$ constant matrix, and $\mathbf{1}_p$ is a p -dimensional vector with all elements being 1.

It is useful to identify several typical members in this class, such as the OLS estimator and the WLS estimator. Ma and Wang (2012) showed that for uncensored data, the OLS is derived by viewing the \mathbf{q}_i 's as covariates and $I(T_i \leq t)$ as response variables, where the covariates and the responses are linked by $F(t)$ via a linear regression model

$$Y_i \equiv I(T_i \leq t) = \mathbf{q}_i^T \mathbf{F}(t) + e_i,$$

with $E(e_i | \mathbf{q}_i) = 0$, $i = 1, \dots, n$. It is straightforward that the e_i 's are independent, conditional on \mathbf{q}_i 's, and have variances $v_i = \mathbf{q}_i^T \mathbf{F}(t) \{1 - \mathbf{q}_i^T \mathbf{F}(t)\}$. The WLS is then defined by using the inverse of the variances v_i as weights in a WLS.

Both the OLS and the WLS correspond to special members of the family of all influence functions. Specifically, the OLS has the influence function

$$\phi_{\text{OLS}}(\mathbf{q}, s) = \{E(\mathbf{Q}\mathbf{Q}^T)\}^{-1} \mathbf{q} \{I(s < t) - \mathbf{q}^T \mathbf{F}(t)\},$$

and WLS has the influence function

$$\phi_{\text{WLS}}(\mathbf{q}, s) = \{E(W\mathbf{Q}\mathbf{Q}^T)\}^{-1} w\mathbf{q} \{I(s < t) - \mathbf{q}^T \mathbf{F}(t)\},$$

where W is a random weight variable. Furthermore, by projecting an influence function onto the tangent space, Ma and Wang (2012) derived the EFF corresponding to a semiparametric efficient estimator:

$$\phi_{\text{EFF}}(\mathbf{q}, s) = \frac{\{I(s < t) \mathbf{I}_p - \mathbf{K}\} \mathbf{A}^{-1}(s) \mathbf{q}}{\mathbf{q}^T \mathbf{f}(s)},$$

where $\mathbf{A}(s) = \int \frac{\mathbf{Q}\mathbf{Q}^T p_{\mathbf{Q}}(\mathbf{Q})}{\mathbf{Q}^T \mathbf{f}(s)} d\mu(\mathbf{Q})$, and $\mathbf{K} = \int_{T_1}^{T_2} I(s < t) \mathbf{A}^{-1}(s) ds \{ \int_{T_1}^{T_2} \mathbf{A}^{-1}(s) ds \}^{-1}$.

The form ϕ_{EFF} is known, but may contain unknown nuisance parameters such as the density $\mathbf{f}(\cdot)$. As before, we assume $\mathbf{f}(\cdot)$ is completely unspecified (nonparametric), and thus, is an infinite-dimensional nuisance parameter. Ma and Wang (2012) showed that substituting consistent estimators for the nuisance parameters in ϕ_{EFF} and solving for $\mathbf{F}(t)$ leads to a semiparametric efficient estimator that reaches the semiparametric efficiency bound in the sense of Bickel et al. (1993).

SUPPLEMENTARY MATERIALS

The document provides further technical details and additional numerical results from the simulation studies and application to COHORT study.

[Received June 2011. Revised May 2012.]

REFERENCES

- Arias, E. (2006), "United States Life Tables, 2003," *National Vital Statistics Report*, 54, 1–40. [1333]
- Bang, H., and Tsiatis, A. A. (2000), "Estimating Medical Costs With Censored Data," *Biometrika*, 87, 329–343. [1327]
- (2002), "Median Regression With Censored Cost Data," *Biometrics*, 58, 643–649. [1327]
- Begg, C. B. (2002), "On the Use of Familial Aggregation in Population-Based Case Probands for Calculating Penetrance," *Journal of the National Cancer Institute*, 94, 1221–1226. [1336]
- Bickel, P. J., Klaassen, C. A. J., Ritov, Y., and Wellner, J. A. (1993), *Efficient and Adaptive Estimation for Semiparametric Models*, Baltimore, MD: Johns Hopkins University Press. [1337]
- Chatterjee, N., and Wacholder, S. (2001), "A Marginal Likelihood Approach for Estimating Penetrance From Kin-Cohort Designs," *Biometrics*, 57, 245–252. [1325, 1326]
- Diao, G., and Lin, D. (2005), "Semiparametric Methods for Mapping Quantitative Trait Loci With Censored Data," *Biometrics*, 61, 789–798. [1325]
- Dorsey, E. R., Beck, C., Adams, M., Chadwick, G., de Bleeck, E. A., McCallum, C., Briner, L., Deuel, L., Clarke, A., Stewart, R., Shoulson, I., and the Huntington Study Group (2008), "Communicating Clinical Trial Results to Research Participants," *Archives of Neurology*, 65, 1590–1595. [1325]
- Fine, J., Zhou, F., and Yandell, B. (2004), "Nonparametric Estimation of the Effects of Quantitative Trait Loci," *Biostatistics*, 5, 501–513. [1325]
- Foroud, T., Gray, J., Ivashina, J., and Conneally, P. (1999), "Differences in Duration of Huntingtons Disease Based on Age at Onset," *Journal of Neurology, Neurosurgery, and Psychiatry*, 66, 52–56. [1325]
- Harper, P. S. (1996), *Huntington's Disease* (2nd ed.), London: W.B. Saunders. [1334]
- Horvitz, D. G., and Thompson, D. J. (1952), "A Generalization of Sampling Without Replacement From a Finite Universe," *Journal of the American Statistical Association*, 47, 663–685. [1327]
- Huntington's Disease Collaborative Research Group (1993), "A Novel Gene Containing a Trinucleotide Repeat That Is Expanded and Unstable on Huntingtons Disease Chromosomes," *Cell*, 72, 971–983. [1325]
- Khoury, M., Beaty, H., and Cohen, B. (1993), *Fundamentals of Genetic Epidemiology*, New York: Oxford University Press. [1325, 1326]
- Lander, E. S., and Botstein, D. (1989), "Mapping Mendelian Factors Underlying Quantitative Traits Using RFLP Linkage Maps," *Genetics*, 121, 743–756. [1324]
- Langbehn, D., Brinkman, R., Falush, D., Paulsen, J., and Hayden, M. (2004), "A New Model for Prediction of the Age of Onset and Penetrance for Huntingtons Disease Based on CAG Length," *Clinical Genetics*, 65, 267–277. [1325]
- Li, H., Yang, P., and Schwartz, A. G. (1998), "Analysis of Age of Onset Data From Case-Control Family Studies," *Biometrics*, 54, 1030–1039. [1324]
- Lipsitz, S. R., Ibrahim, J. G., and Zhao, L. P. (1999), "A Weighted Estimating Equation for Missing Covariate Data With Properties Similar to Maximum Likelihood," *Journal of the American Statistical Association*, 94, 1147–1160. [1329]
- Ma, Y., and Wang, Y. (2012), "Efficient Semiparametric Estimation for Mixture Data," *Electronic Journal of Statistics*, 6, 710–737. [1325, 1330, 1336, 1337]
- Mai, P. L., Chatterjee, N., Hartge, P., Tucker, M., Brody, L., Struewing, J. P., and Wacholder, S. (2009), "Potential Excess Mortality in BRCA1/2 Mutation Carriers Beyond Breast, Ovarian, Prostate, and Pancreatic Cancers, and Melanoma," *PLoS ONE*, 4, e4812. [1324]
- Marder, K., Levy, G., Louis, E. D., Mejia-Santana, H., Cote, L., Andrews, H., Harris, J., Waters, C., Ford, B., Frucht, S., Fahn, S., and Ottman, R. (2003), "Accuracy of Family History Data on Parkinson's Disease," *Neurology*, 61, 18–23. [1324]
- Marder, K., Zhao, H., Myers, R., Cudkovic, M., Kayson, E., Kiebertz, K., Orme, C., Paulsen, J., Penney, J., Siemers, E., Shoulson, I., and Group, T. H. S. (2000), "Rate of Functional Decline in Huntingtons Disease," *Neurology*, 369, 452–458. [1334]
- Rabinowitz, D. (2000), "Computing the Efficient Score in Semi-Parametric Problems," *Statistica Sinica*, 10, 265–280. [1328]
- Robins, J. M., and Rotnitzky, A. (1992), "Recovery of Information and Adjustment for Dependent Censoring Using Surrogate Markers," in *AIDS Epidemiology*, eds. N. Jewell, K. Dietz, and V. Farewell, Boston, MA: Birkhäuser, pp. 297–331. [1328, 1331]
- Robins, J. M., Rotnitzky, A., and Zhou, L. P. (1994), "Estimation of Regression Coefficients When Some Regressors Are Not Always Observed," *Journal of the American Statistical Association*, 89, 846–866. [1327, 1328]
- Robins, J. M., and Wang, N. (2000), "Inference for Imputation Estimators," *Biometrika*, 87, 113–124. [1329]
- Rubinsztein, D. C., Leggo, J., Coles, R., Almqvist, E., Biancalana, V., Cassiman, J. J., Chotai, K., Connarty, M., Crauford, D., Curtis, A., Curtis, D., et al. (1996), "Phenotypic Characterization of Individuals With 30–40 CAG Repeats in the Huntington Disease (HD) Gene Reveals HD Cases With 36 Repeats and Apparently Normal Elderly Individuals With 36–39 Repeats," *American Journal of Human Genetics*, 59, 16–22. [1325]
- Struewing, J. P., Hartge, P., Wacholder, S., Baker, S. M., Berlin, M., McAdams, M., Timmerman, M. M., Brody, L. C., and Tucker, M. A. (1997), "The

- Risk of Cancer Associated With Specific Mutations of BRCA1 and BRCA2 Among Ashkenazi Jews," *The New England Journal of Medicine*, 336, 1401–1408. [1324]
- Tsiatis, A. A., and Ma, Y. (2004), "Locally Efficient Semiparametric Estimators for Functional Measurement Error Model," *Biometrika*, 91, 835–848. [1328]
- Van der Laan, M. J., and Hubbard, A. E. (1998), "Locally Efficient Estimation of the Survival Distribution With Right-Censored Data and Covariates When Collection of Data Is Delayed," *Biometrika*, 85, 771–783. [1328]
- Wacholder, S., Hartge, P., Struwing, J., Pee, D., McAdams, M., Brody, L., and Tucker, M. (1998), "The Kin-Cohort Study for Estimating Penetrance," *American Journal of Epidemiology*, 148, 623–630. [1324,1325,1326,1333]
- Wang, N., and Robins, J. M. (1998), "Large Sample Inference in Parametric Multiple Imputation," *Biometrika*, 85, 935–948. [1329]
- Wang, Y., Clark, L. N., Louis, E. D., Mejia-Santana, H., Harris, J., Cote, L. J., Waters, C., Andrews, H., Ford, B., Frucht, S., Fahn, S., Ottman, R., Rabinowitz, D., and Marder, K. (2008), "Risk of Parkinson Disease in Carriers of Parkin Mutations: Estimation Using the Kin-Cohort Method," *Archives of Neurology*, 65, 467–474. [1324,1333]
- Wang, Y., Clark, L. N., Marder, K., and Rabinowitz, D. (2007), "Non-Parametric Estimation of Genotype-Specific Age-at-Onset Distributions From Censored Kin-Cohort Data," *Biometrika*, 94, 403–414. [1325]
- Wu, R., Ma, C., and Casella, G. (2007), *Statistical Genetics of Quantitative Traits: Linkage, Maps, and QTL*, New York: Springer. [1324,1326]
- Wu, R., Ma, C., Chang, M., Littell, R., Wu, S., Yin, T., Huang, M., Wang, M., and Casella, G. (2002), "A Logistic Mixture Model for Characterizing Genetic Determinants Causing Differentiation in Growth Trajectories," *Genetical Research*, 79, 235–245. [1325]
- Wu, R., Zeng, Z., McKend, S., and OMalley, D. (2000), "The Case for Molecular Mapping in Forest Tree Breeding," *Plant Breeding Reviews*, 19, 41–68. [1325]
- Yu, Z., and Lin, X. (2008), "Nonparametric Regression Using Local Kernel Estimating Equations for Correlated Failure Time Data," *Biometrika*, 95, 123–137. [1325]
- Zeng, D., and Lin, D. (2007), "Maximum Likelihood Estimation in Semiparametric Models With Censored Data" (with discussion), *Journal of the Royal Statistical Society, Series B*, 69, 507–564. [1325,1336]
- Zhang, M., Tsiatis, A. A., and Davidian, M. (2008), "Improving Efficiency of Inferences in Randomized Clinical Trials Using Auxiliary Covariates," *Biometrics*, 64, 707–715. [1328]
- Zhao, W., and Wu, R. (2008), "Wavelet-Based Nonparametric Functional Mapping of Longitudinal Curves," *Journal of the American Statistical Association*, 103, 714–725. [1325]