



Published in final edited form as:

*Biometrika*. 2014 June ; 101(2): 465–476. doi:10.1093/biomet/asu004.

## Bootstrap for the case-cohort design

Yijian Huang

Department of Biostatistics and Bioinformatics, Rollins School of Public Health, Emory University, Atlanta, Georgia 30322, U.S.A

Yijian Huang: yhuang5@emory.edu

### Summary

The case-cohort design facilitates economical investigation of risk factors in a large survival study, with covariate data collected only from the cases and a simple random subset of the full cohort. Methods that accommodate the design have been developed for various semiparametric models, but most inference procedures are based on asymptotic distribution theory. Such inference can be cumbersome to derive and implement, and does not permit confidence band construction. While bootstrap is an obvious alternative, how to resample is unclear because of complications from the two-stage sampling design. We establish an equivalent sampling scheme, and propose a novel and versatile nonparametric bootstrap for robust inference with an appealingly simple single-stage resampling. Theoretical justification and numerical assessment are provided for a number of procedures under the proportional hazards model.

### Keywords

Confidence band; Interval estimation; Multiplier bootstrap; Proportional hazards model; Robust inference; Simple random sampling; Two-stage sampling

## 1. Introduction

In clinical or epidemiologic investigation of infrequent disease endpoints, cohort studies require a large sample size, but full-cohort covariate data collection can be costly. Case-cohort sampling (Prentice, 1986) offers an economical alternative, to collect covariates only from the cases and a simple random subcohort. A number of statistical methods have been developed for the proportional hazards model under such sampling. Prentice (1986) and Self & Prentice (1988) proposed a pseudolikelihood approach. Kalbfleisch & Lawless (1988) and Chen & Lo (1999) suggested methods to take better advantage of the data so as to improve estimation efficiency. Borgan et al. (2000), Kulich & Lin (2004), and Nan (2004) studied a variant of the design where some covariates or surrogate measures are available for the full cohort. Other semiparametric models have also been investigated, including the additive hazards model (Kulich & Lin, 2000), linear transformation models (Chen, 2001; Kong et al., 2004; Lu & Tsiatis, 2006; Chen & Zucker, 2009), and the accelerated failure time model (Nan et al., 2006; Kong & Cai, 2009).

Statistical challenges with the case-cohort design lie in the two-stage sampling. After the first stage resulting in the full study cohort, the second stage selects a subcohort by simple random sampling without replacement. As a result, interval estimation based on asymptotic

distribution theory, which is routinely adopted and specific to each point estimation procedure, is often cumbersome to derive and implement. Although the jackknife approach of Barlow (1994) to variance estimation in the proportional hazards model may deal with some of these issues, neither approach permits confidence band construction for an infinite-dimensional quantity, e.g., a covariate-specific survival function for prediction. While a bootstrap method would be advantageous in these regards, how to resample is not obvious since the covariates are not available for the full cohort. The only existing bootstrap method is due to Wacholder et al. (1989), which has been adopted for the proportional hazards model and accelerated failure time model (Kong & Cai, 2009). However, the method fixes the case numbers in the subcohort and full cohort. Theoretical justification is lacking and it may not always perform well.

We propose a novel nonparametric bootstrap that involves only a single-stage resampling, after establishing an equivalent sampling scheme for the case-cohort design. Three existing point estimation methods under the proportional hazards model will set the stage and illustrate our proposal.

## 2. Point estimation methods under the proportional hazards model

Write the survival time as  $T$  and the censoring time as  $C$ . As a result of censoring, they are observed only through follow-up time  $X = T \wedge C$  and censoring indicator  $\delta = I(T \leq C)$ , where  $\wedge$  is the minimization operator and  $I(\cdot)$  is the indicator function. Denote the covariate vector by  $Z$ ; we confine our attention to time-independent covariates. The proportional hazards model (Cox, 1972) postulates

$$d\Lambda_z(t) = e^{\beta_0^T z} d\Lambda_0(t), \quad T \perp\!\!\!\perp C | Z, \quad (1)$$

where  $\Lambda_z$  is the cumulative hazard function of  $T$  given  $Z = z$ ,  $\Lambda_0$  is an unspecified baseline cumulative hazard function,  $\beta_0$  is an unknown regression coefficient, and  $\perp\!\!\!\perp$  denotes statistical independence. Under cohort sampling, the data consist of  $(X_i, \delta_i, Z_i), i = 1, \dots, n$ , as  $n$  independent replicates of  $(X, \delta, Z)$ . Define counting process  $N_i(t) = \delta_i I(X_i \leq t)$  and at-risk process  $Y_i(t) = I(X_i > t)$ . Further, introduce empirical processes

$$\begin{aligned} A(t) &= \frac{1}{n} \sum_{i \in \mathcal{F}} N_i(t), & B(t) &= \frac{1}{n} \sum_{i \in \mathcal{F}} Z_i N_i(t), \\ U(t, \beta) &= \frac{1}{n} \sum_{i \in \mathcal{F}} Y_i(t) e^{\beta^T Z_i}, & V(t, \beta) &= \frac{1}{n} \sum_{i \in \mathcal{F}} Y_i(t) Z_i e^{\beta^T Z_i}, \end{aligned}$$

where  $\mathcal{F} \equiv \{1, \dots, n\}$  is the index set. The maximum partial likelihood estimator (Cox, 1972, 1975),  $\hat{\beta}$ , is the solution of

$$\Psi(\beta) = B(\infty) - \int_0^\infty \frac{V(s, \beta)}{U(s, \beta)} dA(s).$$

The Breslow (1972) estimator of  $\Lambda_0(t)$  is  $\hat{\Lambda}(t; \hat{\beta})$ , with

$$\hat{\Lambda}(t; \hat{\beta}) = \int_0^t \frac{dA(s)}{U(s, \hat{\beta})}$$

These functional representations are due to Huang & Wang (2000).

However, under the case-cohort design, the covariates are ascertained only for the cases, for whom  $i = 1$ , and a simple random subcohort  $\mathcal{S} \subset \mathcal{F}$  of size  $m$ . Thus,  $U(t, \hat{\beta})$  and  $V(t, \hat{\beta})$  in  $\Psi(\hat{\beta})$  and  $\hat{\Lambda}(t; \hat{\beta})$  are no longer available; they are the full-cohort empirical estimates of  $E\{Y(t)e^{\beta^T Z}\}$  and  $E\{Y(t)Ze^{\beta^T Z}\}$ , respectively. Different replacements are made in the following three case-cohort estimation procedures, giving rise to

$$\Psi_k(\hat{\beta}) = B(\infty) - \int_0^\infty \frac{V_k(s, \hat{\beta})}{U_k(s, \hat{\beta})} dA(s), \quad \hat{\Lambda}_k(t; \hat{\beta}) = \int_0^t \frac{dA(s)}{U_k(s, \hat{\beta})}, \quad (2)$$

and subsequently estimator  $\{\hat{\beta}_k, \hat{\Lambda}_k(t; \hat{\beta}_k)\}$ , for  $k = 1, 2, 3$ . Self & Prentice (1988) adopted the subcohort counterparts,

$$U_1(t, \hat{\beta}) = \frac{1}{m} \sum_{i \in \mathcal{S}} Y_i(t) e^{\beta^T Z_i}, \quad V_1(t, \hat{\beta}) = \frac{1}{m} \sum_{i \in \mathcal{S}} Y_i(t) Z_i e^{\beta^T Z_i}.$$

Chen & Lo (1999) showed that one can make better use of the case-cohort data. Both  $E\{Y(s)e^{\beta^T Z}\}$  and  $E\{Y(s)Ze^{\beta^T Z}\}$  are weighted averages of case- and control-specific quantities, e.g.,  $E\{Y(s)e^{\beta^T Z}\} = E(\cdot)E\{Y(s)e^{\beta^T Z} | i = 1\} + \{1 - E(\cdot)\}E\{Y(s)e^{\beta^T Z} | i = 0\}$ , and estimating the case-specific quantities may utilize the full cohort instead of the subcohort. By taking  $m_1/m$  as an estimate of  $E(\cdot)$  with  $m_1 = \sum_{i \in \mathcal{S}} i$ , their first method adopts

$$U_2(t, \hat{\beta}) = \frac{m_1}{m} \frac{1}{n_1} \sum_{i \in \mathcal{F}^1} Y_i(t) e^{\beta^T Z_i} + \frac{1}{m} \sum_{i \in \mathcal{S}^0} Y_i(t) e^{\beta^T Z_i}, \\ V_2(t, \hat{\beta}) = \frac{m_1}{m} \frac{1}{n_1} \sum_{i \in \mathcal{F}^1} Y_i(t) Z_i e^{\beta^T Z_i} + \frac{1}{m} \sum_{i \in \mathcal{S}^0} Y_i(t) Z_i e^{\beta^T Z_i},$$

where  $n_1 = \sum_{i \in \mathcal{F}} i$ ,  $\mathcal{F}^1 = \{i \in \mathcal{F} : i = 1\}$ , and  $\mathcal{S}^0 = \{i \in \mathcal{S} : i = 0\}$ . If  $n_1/n$  instead is used as an estimate of  $E(\cdot)$ , their second method has

$$U_3(t, \hat{\beta}) = \frac{1}{n} \sum_{i \in \mathcal{F}^1} Y_i(t) e^{\beta^T Z_i} + \frac{n - n_1}{n} \frac{1}{m - m_1} \sum_{i \in \mathcal{S}^0} Y_i(t) e^{\beta^T Z_i}, \\ V_3(t, \hat{\beta}) = \frac{1}{n} \sum_{i \in \mathcal{F}^1} Y_i(t) Z_i e^{\beta^T Z_i} + \frac{n - n_1}{n} \frac{1}{m - m_1} \sum_{i \in \mathcal{S}^0} Y_i(t) Z_i e^{\beta^T Z_i}.$$

Rescaled estimating functions  $n\Psi_1(\hat{\beta})$  and  $n\Psi_2(\hat{\beta})$ , and therefore  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , do not require the full cohort size,  $n$ . As pointed out in Prentice (1986), the case-cohort design does not

necessarily require a full cohort roster and thus  $n$  need not be known; see also Chen & Lo (1999, Remark 2). However,  $\hat{\beta}_3$  and the three estimators of the baseline cumulative hazard function require that the full cohort be well-defined and  $n$  known.

These regression coefficient estimators, as commonly adopted, have been well studied, and their asymptotics-based inference procedures have been developed (Self & Prentice, 1988; Chen & Lo, 1999). However, inference for the baseline cumulative hazard function or a covariate-specific survival function is only available in a pointwise fashion for the method of Self & Prentice (1988). A nonparametric bootstrap is desirable to permit simple and automatic inference, for these as well as other procedures.

### 3. Equivalent sampling scheme and the proposed bootstrap

Efron's (1979) bootstrap would mimic the two-stage sampling to resample the full cohort as a pseudo-population, but the full cohort is not fully observed and possibly not even well-defined. Therefore, the procedure is not applicable, as recognized by Wacholder et al. (1989). In the sample survey literature, Gross (1980), Bickel & Freedman (1984), Chao & Lo (1985), and Sitter (1992a,b) developed methods to construct a pseudo-population for simple random sampling without replacement. Although these methods can be adapted, the resulting bootstraps may not be ideal, for several reasons. First, the resampling is complex, especially when  $n/m$  is not an integer. Second, cases outside the subcohort are not utilized. Third, this approach does not apply when the full cohort size  $n$  is unknown. Finally, a resample may contain only censored observations. Cohort sampling might suffer this complication as well (e.g., Kosorok, 2008) but it can be particularly acute with typical case-cohort studies, where the endpoint is infrequent and the subcohort has limited size.

Appealing to finite population sampling theory seems natural to deal with simple random sampling without replacement; this tactic is also commonly taken for asymptotic studies (e.g., Chen & Lo, 1999; Kulich & Lin, 2000; Kong et al., 2004). However, the full cohort is a random sample, not a finite population of interest. We rather pursue a direct approach by establishing an equivalent sampling scheme.

#### Proposition 1

The joint distribution of a set of random variables that are independent and identically distributed is invariant to reordering by a random permutation.

Since simple random sampling can be implemented via permutation, the subcohort in the case-cohort design consists of independent and identically distributed random variables, and so does its complement. Furthermore, the two sets are independent of each other. This fact does not contradict the well-known dependence structure from simple random sampling, which is conditional on the full cohort. Write the complement of  $\mathcal{S}$  as  $\bar{\mathcal{S}} = \mathcal{F} \setminus \mathcal{S}$ . Then,  $\{(X_i, Z_i) : i \in \mathcal{S}\}$  are  $m$  independent replicates of  $(X, Z)$ , and  $\{(X_i, Z_i) : i \in \bar{\mathcal{S}}\}$  are  $n - m$  independent replicates of  $(X, Z)$ . This results in an equivalent single-stage parallel sampling scheme.

This sampling equivalence first facilitates a model-free asymptotic study for the three estimation methods introduced in § 2.

### Proposition 2

Re-define  $\beta_0$  as the solution of

$$E\{ZN(\tau)\} - \int_0^\tau \frac{E\{Y(s)Ze^{\beta_0^T Z}\}}{E\{Y(s)e^{\beta_0^T Z}\}} dE\{N(s)\} = 0 \quad (3)$$

and subsequently

$$\Lambda_0(t) = \int_0^t \frac{dE\{N(s)\}}{E\{Y(s)e^{\beta_0^T Z}\}}, \quad t \in [0, \tau],$$

where  $\tau = \sup\{t : \text{pr}(X > t) > 0\}$ . Suppose that the subcohort fraction  $m/n$  converges to a constant  $\rho \in (0, 1)$  as both  $m$  and  $n - m$  approach  $\infty$ , and that the conditions in the Appendix hold. Then, for each  $k = 1, 2, 3$ ,  $\hat{\beta}_k$  is consistent for  $\beta_0$  and  $\hat{\Lambda}_k(t, \hat{\beta}_k)$  is consistent for  $\Lambda_0(t)$  uniformly in  $t \in [0, \tau]$ . In addition,  $n^{1/2}\{\hat{\beta}_k - \beta_0, \hat{\Lambda}_k(\cdot; \hat{\beta}_k) - \Lambda_0(\cdot)\}$  converges weakly to a Gaussian process.

**Remark 1**—This result is slightly more general than those of Self & Prentice (1988) and Chen & Lo (1999) as obtained under the proportional hazards model (1). The model (1) implies the above definitions of  $\beta_0$  and  $\Lambda_0(\cdot)$ , but not vice versa.

More importantly, our proposal of parallel bootstrapping  $\mathcal{S}$  and  $\bar{\mathcal{S}}$  naturally follows. We adapt the multiplier or weighted bootstrap, which assigns a nonnegative random weight to each individual and thus averts the complication of a resample without uncensored observations (cf. Kosorok, 2008). These independent and identically distributed weights,  $\xi_i$  for  $i \in \mathcal{F}$ , are independent of the data and have unit mean and unit variance; the standard exponential distribution was used in all our numerical studies reported later. However, a typical multiplier bootstrap as applied to a single sample standardizes the weights by their average such that the sum is fixed to the sample size (e.g., Kosorok et al., 2004; Kosorok, 2008), leading to the Bayesian bootstrap of Rubin (1981) if the standard exponential distribution is chosen for  $\xi_i$ . In contrast, we do not carry out the standardization and consequently our bootstrap resamples have random sizes  $m^*$  and  $n^* - m^*$  for  $\mathcal{S}$  and  $\bar{\mathcal{S}}$ , respectively, where  $m^* = \sum_{i \in \mathcal{S}} \xi_i$  and  $n^* = \sum_{i \in \mathcal{F}} \xi_i$ . While superfluous in the single-sample case, this modification is critical in the case-cohort design particularly when the full cohort size  $n$  is unknown and thus  $\bar{\mathcal{S}}$  is not well defined. In this circumstance, our bootstrap remains applicable provided that the point estimator is defined.

We now detail the proposed bootstrap for the three estimation methods. Define the bootstrap counterparts,

$$\begin{aligned}
 m_1^* &= \sum_{i \in \mathcal{I}} \xi_i \Delta_i, & n_1^* &= \sum_{i \in \mathcal{I}} \xi_i \Delta_i, \\
 A^*(t) &= \frac{1}{n^*} \sum_{i \in \mathcal{I}} \xi_i N_i(t), & B^*(t) &= \frac{1}{n^*} \sum_{i \in \mathcal{I}} \xi_i Z_i N_i(t), \\
 U_1^*(t, \beta) &= \frac{1}{m^*} \sum_{i \in \mathcal{I}} \xi_i Y_i(t) e^{\beta^T Z_i}, & V_1^*(t, \beta) &= \frac{1}{m^*} \sum_{i \in \mathcal{I}} \xi_i Y_i(t) Z_i e^{\beta^T Z_i}, \\
 U_2^*(t, \beta) &= \frac{m_1^*}{m^* n_1^*} \sum_{i \in \mathcal{I}^1} \xi_i Y_i(t) e^{\beta^T Z_i} + \frac{1}{m^*} \sum_{i \in \mathcal{I}^0} \xi_i Y_i(t) e^{\beta^T Z_i}, \\
 V_2^*(t, \beta) &= \frac{m_1^*}{m^* n_1^*} \sum_{i \in \mathcal{I}^1} \xi_i Y_i(t) Z_i e^{\beta^T Z_i} + \frac{1}{m^*} \sum_{i \in \mathcal{I}^0} \xi_i Y_i(t) Z_i e^{\beta^T Z_i}, \\
 U_3^*(t, \beta) &= \frac{1}{n^*} \sum_{i \in \mathcal{I}^1} \xi_i Y_i(t) e^{\beta^T Z_i} + \frac{n^* - n_1^*}{n^*} \frac{1}{m^* - m_1^*} \sum_{i \in \mathcal{I}^0} \xi_i Y_i(t) e^{\beta^T Z_i}, \\
 V_3^*(t, \beta) &= \frac{1}{n^*} \sum_{i \in \mathcal{I}^1} \xi_i Y_i(t) Z_i e^{\beta^T Z_i} + \frac{n^* - n_1^*}{n^*} \frac{1}{m^* - m_1^*} \sum_{i \in \mathcal{I}^0} \xi_i Y_i(t) Z_i e^{\beta^T Z_i}.
 \end{aligned}$$

The bootstrap estimator  $\{\hat{\beta}_k^*, \hat{\Lambda}_k^*(t; \hat{\beta}_k^*)\}$  results from

$$\Psi_k^*(\beta) = B^*(\infty) - \int_0^\infty \frac{V_k^*(s, \beta)}{U_k^*(s, \beta)} dA^*(s), \quad \hat{\Lambda}_k^*(t; \beta) = \int_0^t \frac{dA^*(s)}{U_k^*(s, \beta)}, \quad k=1, 2, 3. \quad (4)$$

Just like  $\hat{\beta}_1$  and  $\hat{\beta}_2$ , their bootstrap counterparts  $\hat{\beta}_1^*$  and  $\hat{\beta}_2^*$  do not require the knowledge of the full cohort size  $n$ .

**Theorem 1**

Adopt the definitions and conditions in Proposition 2. Suppose that the nonnegative random variable  $\xi_1$  of unit mean and unit variance satisfies  $\int_0^\infty \text{pr}(\xi_1 > x)^{1/2} dx < \infty$ . Conditionally on the data,  $n^{1/2}\{\hat{\beta}_k^* - \hat{\beta}_k, \hat{\Lambda}_k^*(\cdot, \hat{\beta}_k^*) - \hat{\Lambda}_k(\cdot, \hat{\beta}_k)\}$  has the same asymptotic distribution as  $n^{1/2}\{\hat{\beta}_k - \beta_0, \hat{\Lambda}_k(\cdot, \hat{\beta}_k) - \Lambda_0(\cdot)\}$  for each  $k = 1, 2, 3$ .

The proposed bootstrap gives rise to robust inference, similar to Barlow (1994) but different from model-based inference of Self & Prentice (1988) and Chen & Lo (1999). The distribution of  $\{\hat{\beta}_k^* - \hat{\beta}_k, \hat{\Lambda}_k^*(\cdot, \hat{\beta}_k^*) - \hat{\Lambda}_k(\cdot, \hat{\beta}_k)\}$  can be simulated to approximate that of  $\{\hat{\beta}_k - \beta_0, \hat{\Lambda}_k(\cdot, \hat{\beta}_k) - \Lambda_0(\cdot)\}$ . The implementation requires trivial coding beyond the point estimators.

**Remark 2**—Case-cohort sampling specializes to cohort sampling when  $n = m$ . In this case, our proposal reduces to the multiplier bootstrap, which is different from the model-based resampling method of Lin et al. (1994) for the proportional hazards model. Recently, Cheng & Huang (2010) developed general theory for the bootstrap in semiparametric models under cohort sampling. They considered the proportional hazards model as an example, but focused only on the regression coefficients.

**4. Numerical studies**

We simulated under a proportional hazards model with a constant baseline hazard and two covariates. The two covariates were independent, each with a uniform distribution between  $-1$  and  $1$ ; their coefficients were both unity. The censoring time depended linearly on the

first covariate, having a uniform distribution between 0 and 1.5. The full cohort size was 1000. As a realistic scenario, the baseline hazard was set to 12.5, resulting in approximately 90% censoring, and the subcohort size was set to 200. For a more comprehensive assessment, two additional scenarios were studied as well. One reduced the subcohort size to 100, while the other changed the baseline hazard to 1 for a censoring rate of approximately 50%. With these three scenarios, the expected sizes of combined cases and subcohort were 280, 190, and 600, including 100, 100, and 500 expected cases, respectively.

To compare with the proposed bootstrap, we also evaluated asymptotics-based inference for the regression coefficients  $\beta_1$  and  $\beta_2$  as in Chen & Lo (1999, Remark 5) and the bootstrap of Wacholder et al. (1989); asymptotics-based pointwise inference for the baseline cumulative hazard function  $\Lambda_0(\cdot)$  would require tedious derivation and was not examined. For both bootstrap methods, we computed standard errors, Wald-type and percentile confidence intervals for  $\beta_1$ ,  $\beta_2$ ,  $a_1 = \log \Lambda_0(0.5)$ , and  $a_2 = \log \Lambda_0(1)$ . Moreover, a confidence band for the baseline survival function  $S_0(\cdot)$  over time  $[0, 1.25]$  was constructed, by transformation from the equal-precision band for  $\log \Lambda_0(\cdot)$ ; an equal-precision band has boundaries parallel to those of the corresponding pointwise Wald-type confidence intervals. The calculation was based on 1000 bootstrap resamples, but our numerical experiments indicated that a much smaller bootstrap size, say, 200, would typically suffice for standard errors and Wald-type confidence intervals (cf. Efron & Tibshirani, 1993, § 6.4).

Table 1 reports the results from 2000 replications. The proposed bootstrap performed well overall, more so with a larger subcohort; the performance for  $\beta_1$  and  $\beta_2$  was largely comparable to that of the asymptotics-based inference procedures. The standard errors all tracked the standard deviations closely. The Wald-type and percentile confidence intervals had coverage probabilities reasonably close to the nominal level, but the latter slightly outperformed the former for  $a_1$  and  $a_2$ . The confidence bands for  $S_0(\cdot)$  also had good coverage probabilities. In comparison, the bootstrap of Wacholder et al. was less satisfactory except for  $\beta_1$  and  $\beta_2$  in the circumstance of 90% censoring. At 50% censoring, the second method of Chen & Lo was noticeably more efficient than their first method for  $\beta_1$  estimation. However, the bootstrap of Wacholder et al. showed little difference in the standard errors. In addition, the confidence intervals for  $a_1$  and  $a_2$  often under-covered whereas the confidence bands for  $S_0(\cdot)$  tended to over-cover.

Typical case-cohort studies involve infrequent disease endpoints, and have a fairly large full cohort and a moderate-sized subcohort, say, in thousands and hundreds, respectively. The preceding simulations and our other numerical experience suggest that the proposed bootstrap is generally reliable in such circumstances. Furthermore, the bootstrap also performs well with more frequent disease endpoints.

As an illustration, we analyzed data from the ACTG 175 trial conducted by the AIDS Clinical Trials Group (Hammer et al., 1996). The study evaluated four treatments, zidovudine, zidovudine plus didanosine, zidovudine plus zalcitabine, and didanosine, in HIV-1 infected adults whose screening CD4 counts were between 200 and 500 per cubic millimeter. A total of 2467 participants were randomized, and the mean follow-up was 29 months with 154 deaths observed. We considered a survival model with two continuous

covariates, age and  $\log(\text{CD4})$ , and five binary ones, treatment indicators, hemophilia, and presence of symptomatic HIV infection. All these baseline covariates were measured in the dataset, and the full cohort was sampled to emulate the case-cohort design. We began with fitting the proportional hazards model to the full cohort, and subsequently drew 100 subcohorts of size 240 and averaged the case-cohort point estimates and standard errors. The results are summarized in Table 2, where the bootstrap size of 1000 was used. The case-cohort estimates were, on average, all close to those from the full-cohort analysis, but more variable. The two methods of Chen & Lo had estimates comparable to each other and more efficient than those of Self & Prentice. Consistent with the earlier simulation results, the standard errors from our bootstrap were very similar to the asymptotics-based ones for all three case-cohort methods. We also estimated the survival function for an individual with given covariates and constructed a 95% confidence band by the proposed bootstrap, using the same approach as in the earlier simulations for the baseline survival function. Figure 1 shows the averaged survival functions and averaged confidence bands over the 100 simulated subcohorts, for the three case-cohort methods. The confidence bands for the two Chen & Lo methods were barely distinguishable, and tighter than that for Self & Prentice.

## 5. Discussion

Despite using simple random sampling without replacement in the second stage, the case-cohort design gives rise to an independent data structure. This result facilitates statistical developments using standard tools such as empirical process theory. In the literature, Bernoulli sampling has been suggested in place of simple random sampling, partly to have an independent and identically distributed sample so as to exploit standard theory (e.g., Kulich & Lin, 2000, 2004; Nan et al., 2006). It is now clear that this alteration may be unnecessary for the purpose.

The three procedures in § 2 are among a large collection of case-cohort estimation methods for the proportional hazards model. Often, one method does not dominate another in both feasibility and efficiency. For example, the first method of Chen & Lo is more efficient than Self & Prentice when estimating the regression coefficients. However, with time-dependent covariates, the former requires each case outside the subcohort to have its whole covariate history available. In contrast, the latter only needs the covariate at the failure time, which is more realistic particularly with prospective sampling. Another reason for the co-existence of many methods is to accommodate sampling variations, e.g., stratified sampling (Borgan et al., 2000). A general and automatic tool for inference, such as the proposed bootstrap, is thus particularly attractive.

We have focused on the proportional hazards model because of its popularity. The proposed bootstrap should apply to other models with justifications similar to Theorem 1. Under the framework of our modularized proof as given, essentially it suffices to establish that a new estimator is a well-behaved and sufficiently smooth functional of empirical processes. This is clearly the case for the estimators of Kulich & Lin (2000) under the additive hazards model. However, those of Lu & Tsiatis (2006) under linear transformation models and Nan et al. (2006) and Kong & Cai (2009) under the accelerated failure time model may challenge the proof. With the former, an explicit profile estimating function for the finite-dimensional



parameter may not exist. In the accelerated failure time model case, the estimating functions are not smooth. Nevertheless, existing techniques may be adopted or adapted to address these complications.

## Acknowledgments

The author thanks Professors Victor DeGruttola and Michael Hughes for permission to use the ACTG 175 trial data, and the reviewers and Professor Brent Johnson for helpful comments and suggestions that have led to an improved exposition. Partial support by grants from the U.S. National Science Foundation and National Institutes of Health is gratefully acknowledged.

## References

- Barlow WE. Robust variance estimation for the case-cohort design. *Biometrics*. 1994; 50:1064–72. [PubMed: 7786988]
- Bickel PJ, Freedman DA. Asymptotic normality and the bootstrap in stratified sampling. *Ann Statist*. 1984; 12:470–82.
- Borgan Ø, Langholz B, Samuelsen SO, Goldstein L, Pogoda J. Exposure stratified case-cohort designs. *Lifetime Data Anal*. 2000; 6:39–58. [PubMed: 10763560]
- Breslow NE. Discussion of the paper by D. R. Cox. *JR Statist Soc B*. 1972; 34:216–7.
- Cox DR. Regression models and life tables (with Discussion). *J R Statist Soc B*. 1972; 34:187–220.
- Cox DR. Partial likelihood. *Biometrika*. 1975; 62:269–76.
- Chao MT, Lo SH. A bootstrap method for finite populations. *Sankhya A*. 1985; 47:399–405.
- Chen HY. Weighted semiparametric likelihood method for fitting a proportional odds regression model to data from the case-cohort design. *J Am Statist Assoc*. 2001; 96:1446–57.
- Chen K, Lo SH. Case-cohort and case-control analysis with Cox's model. *Biometrika*. 1999; 86:755–64.
- Chen YH, Zucker DM. Case-cohort analysis with semiparametric transformation models. *J Statist Plann Inference*. 2009; 139:3706–17.
- Cheng G, Huang JZ. Bootstrap consistency for general semiparametric  $M$ -estimation. *Ann Statist*. 2010; 38:2884–915.
- Efron B. Bootstrap methods: Another look at the jackknife. *Ann Statist*. 1979; 7:1–26.
- Efron, B.; Tibshirani, RJ. *An Introduction to the Bootstrap*. New York: Chapman & Hall; 1993.
- Gill RD. Non- and semi-parametric maximum likelihood estimators and the von Mises method — I. *Scand J Statist*. 1989; 16:97–128.
- Gross, S. Proc Section on Survey Research Methods. American Statistical Association; Alexandria, VA: 1980. Median estimation in sample surveys; p. 181-4.
- Hammer SM, Katzenstein DA, Hughes MD, Gundacker H, Schooley RT, Haubrich RH, Henry WK, Lederman MM, Phair JP, Niu M, Hirsch MS, Merigan TC. 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*. 1996; 335:1081–90. [PubMed: 8813038]
- Huang Y, Wang CY. Cox regression with accurate covariates unascertainable: A nonparametric-correction approach. *J Am Statist Assoc*. 2000; 95:1209–19.
- Kalbfleisch JD, Lawless JF. Likelihood analysis of multi-state models for disease incidence and mortality. *Statist Med*. 1988; 7:149–60.
- Kong L, Cai J. Case-cohort analysis with accelerated failure time model. *Biometrics*. 2009; 65:135–42. [PubMed: 18537948]
- Kong L, Cai J, Sen PK. Weighted estimating equations for semiparametric transformation models with censored data from a case-cohort design. *Biometrika*. 2004; 91:305–19.
- Kosorok, MR. *Introduction to Empirical Processes and Semiparametric Inference*. New York: Springer; 2008.

- Kosorok MR, Lee BL, Fine JP. Robust inference for univariate proportional hazards frailty regression models. *Ann Statist.* 2004; 32:1448–91.
- Kulich M, Lin DY. Additive hazards regression for case-cohort studies. *Biometrika.* 2000; 87:73–87.
- Kulich M, Lin DY. Improving the efficiency of relative-risk estimation in case-cohort studies. *J Am Statist Assoc.* 2004; 99:832–44.
- Lin DY, Fleming TR, Wei LJ. Confidence bands for survival curves under the proportional hazards model. *Biometrika.* 1994; 81:73–81.
- Lu W, Tsiatis AA. Semiparametric transformation models for the case-cohort study. *Biometrika.* 2006; 93:207–14.
- Nan B. Efficient estimation for case-cohort studies. *Canad J Statist.* 2004; 32:403–19.
- Nan B, Yu M, Kalbfleisch JD. Censored linear regression for case-cohort studies. *Biometrika.* 2006; 93:747–62.
- Prentice RL. A case-cohort design for epidemiologic cohort studies and disease prevention trials. *Biometrika.* 1986; 73:1–11.
- Rubin DB. The Bayesian bootstrap. *Ann Statist.* 1981; 9:130–4.
- Self SG, Prentice RL. Asymptotic distribution theory and efficiency results for case-cohort studies. *Ann Statist.* 1988; 16:64–81.
- Sitter RR. A resampling procedure for complex survey data. *J Am Statist Assoc.* 1992a; 87:755–65.
- Sitter RR. Comparing three bootstrap methods for survey data. *Canad J Statist.* 1992b; 20:135–54.
- Wacholder S, Gail MH, Pee D, Brookmeyer R. Alternative variance and efficiency calculations for the case-cohort design. *Biometrika.* 1989; 76:117–23.

## Appendix

### Proofs of Proposition 2 and Theorem 1

We impose the following fairly standard regularity conditions:

#### Condition A1

The upper support point  $\tau$  of  $X$  is finite. Further,  $\text{pr}(T > \tau) > 0$  and  $\text{pr}(C = \tau) > 0$ .

#### Condition A2

The covariate  $Z$  is bounded.

#### Condition A3

The matrix

$$D = \int_0^\tau \left( \frac{E\{Y(s)Z^{\otimes 2}e^{\beta_0^T Z}\}}{E\{Y(s)e^{\beta_0^T Z}\}} - \left[ \frac{E\{Y(s)Ze^{\beta_0^T Z}\}}{E\{Y(s)e^{\beta_0^T Z}\}} \right]^{\otimes 2} \right) dE\{N(s)\}$$

is nonsingular, where  $v^{\otimes 2} \equiv vv^T$  for vector  $v$ .

Condition A1 is adopted to avoid lengthy technical tail treatment.

By Proposition 1, the distribution of  $\{\hat{\beta}_k, \hat{\Lambda}_k(\cdot, \hat{\beta}_k), \hat{\beta}_k^*, \hat{\Lambda}_k^*(\cdot, \hat{\beta}_k^*)\}$  is the same under either case-cohort or single-stage parallel sampling, for any  $k = 1, 2, 3$ . Thus, it suffices to prove Proposition 2 and Theorem 1 under the latter sampling scheme, as we do below. We only

present the proofs for the case of  $k = 2$ , the first method of Chen & Lo (1999). Those for the other two methods are similar and thus omitted.

**Proof of Proposition 2**

We express  $\Psi_2(\beta)$  and  $\Lambda_2(t; \beta)$  as functionals of empirical processes, and exploit empirical process theory. Such an approach was taken by Huang & Wang (2000) and Kosorok (2008, § 4.2.1) under cohort sampling, and becomes feasible and effective for the case-cohort design under its equivalent single-stage parallel sampling scheme. The empirical processes are defined for either the sub-cohort  $\mathcal{S}$  or its complement  $\bar{\mathcal{S}}$ . Specifically, the four processes in  $\Psi_2(\beta)$  and  $\Lambda_2(t; \beta)$  as given by (2) can be written as

$$\begin{aligned}
 A(t) &= \frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} N_i(t) + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} N_i(t), \\
 B(t) &= \frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} Z_i N_i(t) + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} Z_i N_i(t), \\
 U_2(t, \beta) &= \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i\right) \frac{\frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i Y_i(t) e^{\beta^T Z_i} + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} \Delta_i Y_i(t) e^{\beta^T Z_i}}{\frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} \Delta_i} + \frac{1}{m} \sum_{i \in \mathcal{S}} (1 - \Delta_i) Y_i(t) e^{\beta^T Z_i}, \\
 V_2(t, \beta) &= \left(\frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i\right) \frac{\frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i Y_i(t) Z_i e^{\beta^T Z_i} + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} \Delta_i Y_i(t) Z_i e^{\beta^T Z_i}}{\frac{m}{n} \frac{1}{m} \sum_{i \in \mathcal{S}} \Delta_i + \left(1 - \frac{m}{n}\right) \frac{1}{n-m} \sum_{i \in \bar{\mathcal{S}}} \Delta_i} + \frac{1}{m} \sum_{i \in \mathcal{S}} (1 - \Delta_i) Y_i(t) Z_i e^{\beta^T Z_i}.
 \end{aligned}$$

Condition A1 effectively limits the time scale to finite interval  $[0, \tau]$ . Let  $\mathcal{B}$  be an arbitrary compact neighborhood of  $\beta_0$ . Under Condition A2, the classes of functions associated with these empirical processes,  $\{N(t) : t \in [0, \tau]\}$ ,  $\{ZN(t) : t \in [0, \tau]\}$ ,  $\{Y(t)e^{\beta^T Z} : t \in [0, \tau], \beta \in \mathcal{B}\}$ ,  $\{(1 - \Delta)Y(t)e^{\beta^T Z} : t \in [0, \tau], \beta \in \mathcal{B}\}$ ,  $\{Y(t)Ze^{\beta^T Z} : t \in [0, \tau], \beta \in \mathcal{B}\}$ , and  $\{(1 - \Delta)Y(t)Ze^{\beta^T Z} : t \in [0, \tau], \beta \in \mathcal{B}\}$ , are all Donsker; see, e.g., Kosorok (2008, § 4.2.1).

Since Donsker implies Glivenko–Cantelli, the empirical processes in  $A(t)$ ,  $B(t)$ ,  $U_2(t, \beta)$ , and  $V_2(t, \beta)$ , converge in probability to their respective limits, uniformly in  $t \in [0, \tau]$  and  $\beta \in \mathcal{B}$  if applicable. By Condition A1, the limit of  $U_2(t, \beta)$  is bounded away from 0 for  $t \in [0, \tau]$  and  $\beta \in \mathcal{B}$ . One can then show that uniformly  $\Psi_2(\beta)$  converges in probability to the left-hand side of (3), which is a monotone function and admits a unique solution  $\beta_0$  by Condition A3. Therefore,  $\hat{\beta}_2$  converges in probability to  $\beta_0$ . The same technique can be used to prove the uniform convergence of  $\Lambda_2(t; \beta)$ . This, coupled with the consistency of  $\hat{\beta}_2$ , establishes that  $\Lambda_2(t; \hat{\beta}_2)$  converges in probability to  $\Lambda_0(t)$  uniformly in  $t \in [0, \tau]$ .

By Taylor expansion,

$$0 = \Psi_2(\hat{\beta}_2) = \Psi_2(\beta_0) + \Psi_2'(\beta_0)(\hat{\beta}_2 - \beta_0) + o_p(\hat{\beta}_2 - \beta_0)$$

since  $\Psi_2''(\beta)$  is bounded by Condition A2. Using similar techniques as before, one can show that  $\Psi_2'(\beta_0)$  converges in probability to  $-D$ , which is nonsingular by Condition A3. Therefore,

$$\hat{\beta}_2 - \beta_0 = D^{-1} \Psi_2(\beta_0) \{1 + o_p(1)\}.$$

Similarly,

$$\begin{aligned} \hat{\Lambda}_2(t; \hat{\beta}_2) - \Lambda_0(t) &= \hat{\Lambda}_2(t; \beta_0) - \Lambda_0(t) + J(t)(\hat{\beta}_2 - \beta_0) + o_p(\hat{\beta}_2 - \beta_0) \\ &= \hat{\Lambda}_2(t; \beta_0) - \Lambda_0(t) + J(t)D^{-1}\Psi_2(\beta_0)\{1 + o_p(1)\}, \end{aligned}$$

where  $J(t)$  is the limit of  $\hat{\Lambda}_2(t; \beta) / \beta_{\beta=\beta_0}$ . Thus, asymptotically  $\{\hat{\beta}_2 - \beta_0, \hat{\Lambda}_2(t; \hat{\beta}_2) - \Lambda_0(t)\}$  is a linear function of  $\{\Psi_2(\beta_0), \hat{\Lambda}_2(t; \beta_0) - \Lambda_0(t)\}$ .

Given the asymptotic linearity, it remains to establish the asymptotic normality of  $\{\Psi_2(\beta_0), \hat{\Lambda}_2(\cdot; \beta_0) - \Lambda_0(\cdot)\}$ . Observe that  $A(\cdot), B(\cdot), U_2(\cdot, \beta_0)$ , and  $V_2(\cdot, \beta_0)$  are asymptotically equivalent when the term  $m/n$  in these expressions is replaced by its limit  $\rho$ . For example,

$$A(t) = \rho \frac{1}{m} \sum_{i \in \mathcal{I}} N_i(t) + (1-\rho) \frac{1}{n-m} \sum_{i \in \mathcal{F}} N_i(t) + \left(\frac{m}{n} - \rho\right) \left\{ \frac{1}{m} \sum_{i \in \mathcal{I}} N_i(t) - \frac{1}{n-m} \sum_{i \in \mathcal{F}} N_i(t) \right\},$$

where the last term is  $o_p(n^{-1/2})$ . The asymptotic normality of  $\{A(\cdot), B(\cdot), U_2(\cdot, \beta_0), V_2(\cdot, \beta_0)\}$  then follows from the properties of Donsker classes and the functional delta method. Using Gill (1989, Lemma 3) and the chain rule, one can show that  $\Psi_2(\beta_0)$  and  $\hat{\Lambda}_2(\cdot; \beta_0)$  are compactly differentiable functionals of  $\{A(\cdot), B(\cdot), U_2(\cdot, \beta_0), V_2(\cdot, \beta_0)\}$ . Finally, the functional delta method leads to the asymptotic normality of  $\{\Psi_2(\beta_0), \hat{\Lambda}_2(\cdot; \beta_0) - \Lambda_0(\cdot)\}$ .

### Proof of Theorem 1

In parallel with  $A(t), B(t), U_2(t, \beta_0)$ , and  $V_2(t, \beta_0)$ , their counterparts  $A^*(t), B^*(t), U_2^*(t, \beta_0)$ , and  $V_2^*(t, \beta_0)$  in (4) involve bootstrapped empirical processes and bootstrapped subcohort fraction  $m^*/n^*$ . By Kosorok (2008, Corollary 10.14), these bootstrapped empirical processes converge in probability to their respective limits, uniformly in  $t \in [0, \tau]$  and  $\beta \in \mathcal{B}$  if applicable. By the law of large numbers and the continuous mapping theorem,  $m^*/n^*$  converges in probability to  $\rho$ . Then, the same argument as in the consistency proof of Proposition 2 establishes the consistency of  $\hat{\beta}_2^*$  for  $\beta_0$  and  $\hat{\Lambda}_2^*(t; \hat{\beta}_2^*)$  for  $\Lambda_0(t)$  uniformly in  $t \in [0, \tau]$ .

Similar to the proof of Proposition 2, we further obtain

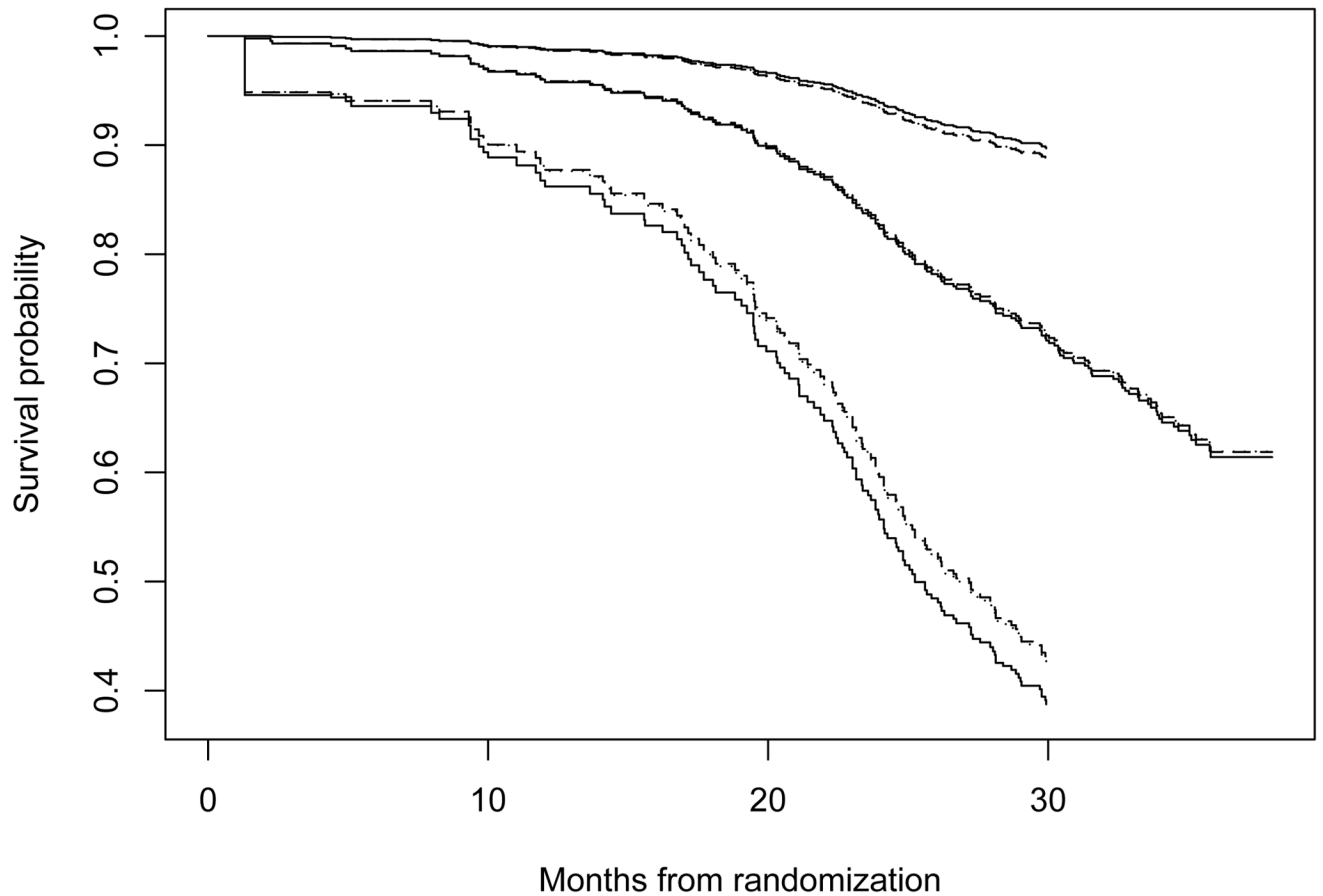
$$\begin{aligned} \hat{\beta}_2^* - \beta_0 &= D^{-1} \Psi_2^*(\beta_0) \{1 + o_p(1)\}, \\ \hat{\Lambda}_2^*(t; \hat{\beta}_2^*) - \Lambda_0(t) &= \hat{\Lambda}_2^*(t; \beta_0) - \Lambda_0(t) + J(t)D^{-1}\Psi_2^*(\beta_0)\{1 + o_p(1)\}. \end{aligned}$$

By Kosorok (2008, Theorem 2.6), the bootstrapped empirical processes in  $A^*(t)$ ,  $B^*(t)$ ,  $U_2^*(t, \beta_0)$ , and  $V_2^*(t, \beta_0)$  minus their expectations are all  $O_p(n^{-1/2})$ . So is  $\Psi_2^*(\beta_0)$ . Then, coupled with the asymptotic linearity result on  $\hat{\beta}_2$  and  $\hat{\Lambda}_2(t; \hat{\beta}_2)$ , we obtain

$$\begin{aligned}\hat{\beta}_2^* - \hat{\beta}_2 &= D^{-1}\{\Psi_2^*(\beta_0) - \Psi_2(\beta_0)\} + o_p(n^{-1/2}), \\ \hat{\Lambda}_2^*(t; \hat{\beta}_2^*) - \hat{\Lambda}_2(t; \hat{\beta}_2) &= \hat{\Lambda}_2^*(t; \beta_0) - \hat{\Lambda}_2(t; \beta_0) + J(t)D^{-1}\{\Psi_2^*(\beta_0) - \Psi_2(\beta_0)\} + o_p(n^{-1/2}).\end{aligned}$$

Thus, asymptotically  $\{\hat{\beta}_2^* - \hat{\beta}_2, \hat{\Lambda}_2^*(t; \hat{\beta}_2^*) - \hat{\Lambda}_2(t; \hat{\beta}_2)\}$  is the same linear function in terms of  $\{\Psi_2^*(\beta_0) - \Psi_2(\beta_0), \hat{\Lambda}_2^*(t; \beta_0) - \hat{\Lambda}_2(t; \beta_0)\}$ , as  $\{\hat{\beta}_2 - \beta_0, \hat{\Lambda}_2(t; \hat{\beta}_2) - \Lambda_0(t)\}$  in terms of  $\{\Psi_2(\beta_0), \hat{\Lambda}_2(t; \beta_0) - \Lambda_0(t)\}$ .

It remains to show that  $n^{1/2}\{\Psi_2^*(\beta_0) - \Psi_2(\beta_0), \hat{\Lambda}_2^*(t; \beta_0) - \hat{\Lambda}_2(t; \beta_0)\}$ , conditionally on the data, has the same asymptotic distribution as  $n^{1/2}\{\Psi_2(\beta_0), \hat{\Lambda}_2(t; \beta_0) - \Lambda_0(t)\}$ . This can be obtained by using a conditional multiplier central limit theorem (Kosorok, 2008, Theorem 2.6) along with the functional delta method and chain rule, as argued in the asymptotic normality proof of Proposition 2.



**Fig. 1.**

Estimated survival functions and 95% confidence bands over  $[0, 30]$  months, as averages over 100 simulated subcohorts, for a 35-year-old individual in the zidovudine arm with symptomatic HIV infection, hemophilia, and baseline CD4 count of 200. The solid, dotted, and dashed lines correspond to the Self & Prentice, the first of Chen & Lo, and the second of Chen & Lo methods, respectively.

Case-cohort simulations to compare the proposed bootstrap (Prop.) with asymptotics-based inference (Asymp.) and the bootstrap of Wacholder et al. (WGPB)

Table 1

	Self & Prentice						Chen & Lo: 1st						Chen & Lo: 2nd					
	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$	$S_0$	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$	$S_0$	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$	
	censoring rate: 90%, subcohort size: 200																	
B	29	31	-43	-29	22	29	-40	-26	19	27	-40	-25						
SD	353	279	208	207	339	269	205	203	332	268	204	201						
SE, Asymp.	350	266			338	257			334	256								
WC, Asymp.	94.7	94.1			94.7	94.0			94.8	94.2								
SE, WGPB	366	282	192	185	347	266	185	176	346	266	184	176						
WC, WGPB	95.9	96.0	93.4	92.1	95.3	95.0	92.5	91.2	95.7	95.1	92.5	91.6						
PC, WGPB	94.9	94.8	90.2	89.8	99.5	94.0	94.2	89.7	99.6	94.2	90.0	89.0	99.5					
SE, Prop.	365	277	219	215	345	261	212	207	338	258	210	205						
WC, Prop.	95.8	95.8	97.0	96.8	95.1	95.0	96.0	96.0	95.3	94.5	95.9	96.0						
PC, Prop.	94.1	93.7	94.6	94.5	96.3	94.4	93.7	94.4	95.8	94.7	93.6	94.4	96.3					
	censoring rate: 90%, subcohort size: 100																	
B	73	60	-50	-41	52	46	-42	-32	53	47	-42	-33						
SD	457	361	233	236	420	328	221	222	409	327	220	219						
SE, Asymp.	436	344			403	319			390	313								
WC, Asymp.	93.8	94.2			93.7	92.9			93.8	93.6								
SE, WGPB	509	410	259	255	417	338	208	202	415	337	208	201						
WC, WGPB	96.9	97.8	96.0	95.2	94.4	95.3	93.8	92.5	94.5	95.3	93.7	93.1						
PC, WGPB	94.5	95.1	92.0	92.0	99.7	93.2	93.8	91.0	99.5	93.6	93.7	90.7	99.7					
SE, Prop.	464	364	252	252	413	325	230	229	397	318	228	224						
WC, Prop.	96.4	97.2	97.2	97.2	94.8	95.0	95.7	96.5	93.7	93.8	95.6	96.3						
PC, Prop.	93.2	93.3	93.7	94.5	96.1	93.2	93.5	94.0	95.8	93.1	93.2	94.2	94.5	95.7				
	censoring rate: 50%, subcohort size: 200																	
B	33	28	-4	13	0	7	-1	8	-1	6	-1	7						
SD	224	205	90	125	151	137	77	104	128	134	76	100						
SE, Asymp.	216	193			149	129			127	126								

	Self & Prentice					Chen & Lo: 1st					Chen & Lo: 2nd				
	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$	$\beta_1$	$\beta_2$	$\alpha_1$	$\alpha_2$	$S_0$
WC, Asymp.	94.2	93.9				94.8	93.8				94.8	92.8			
SE, WGPB	206	199	79	115		122	126	65	87		122	125	65	88	
WC, WGPB	92.8	94.4	91.9	93.0		87.9	92.8	90.3	89.7		93.5	92.8	90.8	91.2	
PC, WGPB	98.9	98.8	95.2	95.8	99.7	87.9	92.8	90.2	88.8	99.4	93.5	93.3	90.5	90.5	99.2
SE, Prop.	219	195	89	126		148	127	76	101		126	123	75	96	
WC, Prop.	95.0	94.0	95.2	95.8		94.5	93.5	94.7	94.5		94.2	92.1	94.5	94.7	
PC, Prop.	93.2	92.5	94.5	93.8	96.8	94.8	93.2	94.7	94.3	94.6	94.2	92.1	94.5	94.0	94.7

B: Empirical bias ( $\times 1000$ ); SD: Empirical standard deviation ( $\times 1000$ ); SE: Average standard error ( $\times 1000$ ); WC: Empirical coverage (%) of 95% Wald-type confidence interval; PC: Empirical coverage (%) of 95% percentile confidence interval or confidence band.

Empty cells indicate either unavailable or inapplicable results.



**Table 2**

Analysis results of the ACTG 175 study

	ZDV+ddI	ZDV+ddC	ddI	age	hemophilia	symptom	log(CD4)
	Full cohort analysis						
estimate	-0.627	-0.391	-0.685	0.039	0.586	0.788	-1.673
SE, Asymp.	0.226	0.209	0.230	0.009	0.273	0.175	0.248
	Case-cohort analyses						
Self & Prentice	-0.699	-0.404	-0.758	0.044	0.740	0.828	-1.841
SE, Asymp.	0.361	0.359	0.361	0.015	0.517	0.310	0.436
SE, Prop.	0.374	0.362	0.377	0.017	0.497	0.320	0.473
Chen & Lo: 1st	-0.692	-0.419	-0.751	0.042	0.678	0.810	-1.811
SE, Asymp.	0.343	0.329	0.345	0.014	0.416	0.279	0.402
SE, Prop.	0.348	0.330	0.350	0.014	0.420	0.285	0.421
Chen & Lo: 2nd	-0.693	-0.419	-0.752	0.043	0.682	0.811	-1.813
SE, Asymp.	0.338	0.324	0.339	0.013	0.409	0.274	0.395
SE, Prop.	0.348	0.330	0.351	0.014	0.416	0.285	0.421

ZDV: zidovudine, ddI: didanosine, ddC: zalcitabine, symptom: presence of symptomatic HIV infection.

SE: standard error. Asympt.: asymptotics-based, Prop.: based on the proposed bootstrap.

Case-cohort estimates and standard errors are averages over 100 simulated subcohorts.