

# Augmented Case-Only Designs for Randomized Clinical Trials with Failure Time Endpoints

James Y. Dai,\* Xinyi Cindy Zhang,\*\* Ching-Yun Wang,\*\*\* and Charles Kooperberg\*\*\*\*

Fred Hutchinson Cancer Research Center and University of Washington, Seattle, Washington

\**email:* jdai@fredhutch.org

\*\**email:* xzhan2@fredhutch.org

\*\*\**email:* cywang@fredhutch.org

\*\*\*\**email:* clk@fredhutch.org

**SUMMARY.** Under suitable assumptions and by exploiting the independence between inherited genetic susceptibility and treatment assignment, the case-only design yields efficient estimates for subgroup treatment effects and gene-treatment interaction in a Cox model. However it cannot provide estimates of the genetic main effect and baseline hazards, that are necessary to compute the absolute disease risk. For two-arm, placebo-controlled trials with rare failure time endpoints, we consider augmenting the case-only design with random samples of controls from both arms, as in the classical case-cohort sampling scheme, or with a random sample of controls from the active treatment arm only. The latter design is motivated by vaccine trials for cost-effective use of resources and specimens so that host genetics and vaccine-induced immune responses can be studied simultaneously in a bigger set of participants. We show that these designs can identify all parameters in a Cox model and that the efficient case-only estimator can be incorporated in a two-step plug-in procedure. Results in simulations and a data example suggest that incorporating case-only estimators in the classical case-cohort design improves the precision of all estimated parameters; sampling controls only in the active treatment arm attains a similar level of efficiency.

**KEY WORDS:** Case-cohort design; Case-only estimator; Gene-treatment interaction; Nested case-control design; Pharmacogenetics.

## 1. Introduction

Individuals respond differently to treatment or prevention modalities, depending on their genetic background, environmental exposures, and clinical characteristics (Charlab and Zhang, 2013). In clinical trials, there is a growing interest to discover and characterize individual or subgroup treatment responses, supplementing primary intent-to-treat analyses. For instance, the emerging pharmacogenetics research aims to identify genetic susceptibility that contributes to inter-individual variability of treatment efficacy and safety, in scales ranging from several candidate genes to the whole genome (Evans and McLeod, 2003; Weinshilboum and Wang, 2004). These studies underscore the potential of personalized medicine, and may also elucidate mechanisms of treatment effect.

To this end, this article pertains to sampling designs for characterizing the influence of pre-treatment biomarkers, e.g., a panel of genetic variants, on treatment effects in randomized clinical trials. Ancillary studies of this nature are increasingly common in the genomic era. However, biomarkers can be expensive to measure. To study the association of biomarkers with relatively uncommon study outcomes, including HIV infection, most cancers, and some cardiovascular events, it is cost-effective to adopt some form of outcome-dependent sampling. Popular outcome-dependent sampling schemes in cohort studies include the nested case-control design

and the case-cohort design (Thomas, 1977; Prentice and Breslow, 1978; Prentice, 1986). Stratified versions of the two sampling designs to oversample certain groups have also been developed for better efficiency (Borgan et al., 2000; Langholz and Borgan, 1995). The properties and utilities of the two designs in cohort studies have been discussed (Self and Prentice, 1988; Langholz and Thomas, 1990).

Consider a two-arm, placebo-controlled randomized prevention trial with a rare failure event. The unique feature is that there is unequivocal design-imposed independence between the treatment assignment and pre-treatment biomarkers, e.g., germline genotypes. Exploiting this independence and assuming censoring being non-informative and independent of randomization arms, case-only methods are more efficient than the two aforementioned designs for estimating gene-treatment interactions and subgroup treatment effects on a rare disease endpoint (Vittinghoff and Bauer, 2006; Dai et al., 2012). These assumptions are better suited for phase III prevention trials where adverse effect is not of concern. Though computed from a logistic model, case-only estimators have the interpretation of hazard ratios in the Cox proportional hazards models. Sensitivity of case-only estimators toward violations of these assumptions has been investigated (Vittinghoff and Bauer, 2006). In recent years, use of case-only methods has started to permeate in prevention trials. See, e.g., trials in the Women's Health Initiative and the HIV

Vaccine Trial Network (Prentice et al., 2010; Dai et al., 2014; Li et al., 2014).

The case-only design, however, does not allow estimation of the full set of parameters in a Cox model. Specifically, neither the genotype main effect nor the cumulative baseline hazard function is estimable from cases alone. These parameters are needed to study the absolute risk of the endpoint for genotype groups in each arm. This limitation hinders interpretation and utility of the estimated gene-treatment interaction, because the estimate of individual absolute risk when treated or when not treated will inform medical counseling and guide treatment selection (Gail et al., 1989; Janes et al., 2011). On the other hand, the traditional case-cohort or nested case-control sampling provides estimates of the baseline hazard and absolute risk, but does not incorporate gene-treatment independence. Leveraging the strengths of both types of designs, we consider augmenting the case-only design to enable estimation of the full set of Cox model parameters. In particular, we focus on variations of the case-cohort design in this article, because it is easy to plan ahead a random subcohort for time-invariant genotypes in clinical trials, and because it has the advantage of accommodating multiple outcomes that may arise in an ancillary study.

Specifically, we consider two scenarios of adding controls to the case-only design, for both of which we can incorporate the case-only estimators in two-step plug-in estimation procedures:

**Scenario I:** Classical case-cohort design with controls drawn from both arms. In essence, this is one way of adding controls to the case-only design. In this scenario, we essentially propose a novel two-step estimation procedure for the classical case-cohort design: the case-only estimator is first used to estimate gene-treatment interaction and treatment main effect, these estimators will then be plugged into established case-cohort estimation methods as offsets. This method allows widely used case-cohort sampling to take advantage of efficient case-only estimators. Our contributions also include an explicit formula of variance estimates for this two-step procedure.

**Scenario II:** Augmented case-only (ACO) design with controls drawn from the active treatment arm only. This is a novel design motivated by vaccine trials, as we will elaborate next. Although primarily driven by scientific rationale, this design is of statistical interest, since only three of the four strata formed by case-control status and randomization arm are sampled. It violates the critical identifiability assumption of non-zero sampling probability for all strata in two-phase sampling (Robins et al., 1994; Breslow et al., 2003). The orthogonality between genotype and randomization arm has to be exploited in order to remedy this anomaly. We show that a similar two-step estimation procedure as for **Scenario I** will identify all parameters in a Cox model, and we show in

the simulations that the estimators are nearly as efficient as those in **Scenario I**.

Scientifically, the motivation for selecting controls only from the active treatment arm (the ACO design in **Scenario II**) comes from studies on host genetics and immune correlates in HIV vaccine trials. It is common to study vaccine-specific immune responses in a pre-specified sample of trial participants in the vaccine arm, as no vaccine-induced immune responses are generated in the placebo arm. Case-cohort sampling is commonly used in this setting (McElrath et al., 2008). Take Li et al. (2014), e.g., if a genotype in the FcγR gene is associated with varying vaccine protection in the RV144 trial, it is useful to investigate whether specific vaccine-induced immune responses are associated with such genotype, in order to understand functionally why the vaccine effect varies by host genetics. Such relationship can only be studied in the vaccine arm. In this sense, concentrating controls in the vaccine arm is cost-effective when a pharmacogenetic study is a component of a systematic approach for understanding treatment effect. Similar rationale applies to high-throughput biomarker studies for better understanding hormone effect in clinical trials in the Women’s Health Initiative (Pitteri et al., 2009).

This article is organized as follows. In Section 2.1 and Section 2.2, we review case-cohort sampling and case-only estimators, respectively. The latter section brings new insights on assumptions required for case-only estimators. In Section 2.3, we show that case-only estimators can be built into a two-step estimation procedure for the case-cohort design. The main parameter of interest we illustrate throughout the article is the genetic main effect. The asymptotic covariance matrix of estimators resulting from the two-step procedure is derived. Extending the results from Section 2.3, we show in Section 2.4 that sampling controls only in the active treatment arm is adequate to estimate all Cox model parameters. For completeness we briefly address the alternative ACO design and the nested case-control sampling in Section 2.5. In Section 3 we compare the efficiency of the proposed designs and estimation methods in simulations, where the standard estimation procedure for a case-cohort design with the same sample size is treated as the benchmark. We present in Section 4 a data example with the standard case-cohort sampling, and we compare standard error estimates resulted from the proposed estimation procedures to the original case-cohort methods. We close with a discussion of the utility of the ACO design and some future work.

## 2. Method

Consider a two-arm, placebo-controlled randomized prevention trial in which participants were followed for evaluating treatment effect on time to certain failure event. Let  $Z$  denote a binary treatment indicator taking the value 1 if the participant is assigned to the active treatment arm, and 0 if assigned to the placebo arm. Let  $G$  denote the baseline biomarker of interest, say an inherited genetic variant, and let  $V$  be a set of pre-treatment variables to be adjusted in risk association. Denote  $Y$  and  $C$  as the failure time and the right-censoring time since randomization, respectively. Given















- marginal genetic association and gene-environment interaction. *American Journal of Epidemiology* **176**, 164–173.
- Evans, W. E. and McLeod, H. L. (2003). Pharmacogenomics- drug disposition, drug targets, and side effects. *The New England Journal of Medicine* **348**, 538–549.
- Gail, M. H., Brinton, L. A., Byar, D. P., Corle, D. K., Green, S. B., Schairer, C., et al. (1989). Projecting individualized probabilities of developing breast cancer for white females who are being examined annually. *Journal of the National Cancer Institute* **81**, 1879–1886.
- Goldstein, L. and Langholz, B. (1992). Asymptotic theory for nested case-control sampling in the cox regression model. *Annals of Statistics* **20**, 1903–1928.
- Han, S. S., Rosenberg, P. S., Garcia-Closas, M., Figueroa, J. D., Silverman, D., Chanock, S. J., et al. (2012). Likelihood ratio test for detecting gene (g)-environment (e) interactions under an additive risk model exploiting g-e independence for case-control data. *American Journal of Epidemiology* **176**, 1060–7.
- Janes, H., Pepe, M. S., Bossuyt, P. M., and Barlow, W. E. (2011). Measuring the performance of markers for guiding treatment decisions. *Annals of Internal Medicine* **154**, 253–259.
- Langholz, B. and Borgan, Y. (1995). Counter-matching: a stratified nested case-control sampling method. *Biometrika* **82**, 69–79.
- Langholz, B. and Thomas, D. C. (1990). Nested case-control and case-cohort methods of sampling from a cohort: a critical comparison. *American Journal of Epidemiology* **131**, 169–176.
- Li, S. S., Gilbert, P. B., Tomaras, G. D., Kijak, G., Ferrari, G., Thomas, R., et al. (2014). Fcgr2c polymorphisms associate with hiv-1 vaccine protection in rv144 trial. *Cancer Epidemiology, Biomarkers & Prevention* **124**, 3879–3890.
- Lin, D. Y. (2000). On fitting cox’s proportional hazards models to survey data. *Biometrika* **87**, 37–47.
- Lin, D. Y. and Wei, L. J. (1989). The robust inference for the cox proportional hazards model. *Journal American Statistical Association* **84**, 1074–1078.
- Lin, D. Y. and Ying, Z. (1993). Cox regression with incomplete covariate measurements. *Journal American Statistical Association* **88**, 1341–1349.
- McElrath, M. J., De Rosa, S. C., Moodie, Z., Dubey, S., Kierstead, L., Janes, H., et al. (2008). Hiv-1 vaccine-induced immunity in the test-of-concept step study: a case-cohort analysis. *Lancet* **372**, 1894–1905.
- Murphy, K. M. and Topel, R. H. (1985). Estimation and inference in two-step econometric models. *Journal of Business & Economic Statistics* **3**, 370–379.
- Nan, B. (2004). Efficient estimation for case-cohort studies. *Canadian Journal of Statistics* **32**, 403–419.
- Newey, W. K. and Powell, J. (1990). Efficient estimation of linear and type i censored regression models under conditional quantile restrictions. *Econometric Theory* **6**, 295–317.
- Pandey, J. P., Namboodiri, A. M., Bu, S., Tapsoba, J. D., Sato, A., and Dai, J. Y. (2013). Immunoglobulin genes and the acquisition of hiv infection in a randomized trail of recombinant adenovirus hiv vaccine. *Virology* **441**, 70–4.
- Pitteri, S. J., Hanash, S. H., Aragaki, A., Amon, L. M., Chen, L., Buson, T. B., et al. (2009). Postmenopausal estrogen and progesterin effects on the serum proteome. *Genome Medicine* **1**(12), 121.
- Prentice, R. L. (1986). A case-cohort design for epidemiologic cohort studies and disease prevention trials. *Biometrika* **73**, 1–11.
- Prentice, R. L. and Breslow, N. E. (1978). Retrospective studies and failure time models. *Biometrika* **65**, 153–158.
- Prentice, R. L., Huang, Y., Hinds, D. A., Peters, U., Cox, D. R., Beilharz, E., et al. (2010). Variation in the fgfr2 gene and the effect of a low-fat dietary pattern on invasive breast cancer. *Cancer Epidemiology, Biomarkers & Prevention* **19**, 74–9.
- Robins, J. M., Rotnitzky, A., and Zhao, L. P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of American Statistical Association* **89**, 846–866.
- Samuelsen, S. O. (1997). A pseudolikelihood approach to analysis of nested case-control studies. *Biometrika* **84**, 379–394.
- Self, S. G. and Prentice, R. L. (1988). Asymptotic distribution theory and efficiency results for case-cohort studies. *Annals of Statistics* **16**, 64–81.
- Tchetgen, E. J. and Robins, J. (2010). The semiparametric case-only estimator. *Biometrics* **66**, 1138–1144.
- Therneau, T. M. and Li, H. (1999). Computing the cox model for case-cohort designs. *Lifetime Data Analysis* **5**, 99–112.
- Thomas, D. C. (1977). Addendum to “methods of cohort analysis: appraisal by application to asbestos mining”. *Journal of Royal Statistical Society, Serial A* **140**, 483–485.
- Vittinghoff, E. and Bauer, D. C. (2006). Case-only analysis of treatment-covariate interactions in clinical trials. *Biometrics* **62**, 769–776.
- Weinshilboum, R. and Wang, L. (2004). Pharmacogenomics: bench to bedside. *Nature Reviews Drug Discovery* **3**, 739–748.

Received November 2014. Revised July 2015. Accepted July 2015.

#### APPENDIX

*The asymptotic linear expansion of the IPW estimator for case-cohort sampling*

For computing the expected covariate values at each event time, those at-risk cases occurring outside of the random sub-cohort can be used with proper sampling weights. Specifically, in estimating function (8), the average term (9) is replaced by

$$\mathbf{S}^{(r)}(\boldsymbol{\beta}_g, T_i) = \frac{1}{n} \sum_{i \in S \cup \mathcal{D}} \frac{1}{\pi_i} R_i(t) \exp(\hat{\boldsymbol{\gamma}}^T \mathbf{X}_{1i} + \boldsymbol{\beta}_g^T \mathbf{X}_{2i}) \mathbf{X}_{2i}^{\otimes r},$$

where

$$\pi_i = \begin{cases} \frac{\sum_{I(\Delta_j=0, j \in \mathcal{S})}}{\sum_{I(\Delta_j=0)}} & \text{if } i \in \mathcal{S} \text{ and } \Delta_i = 0 \\ 1 & \text{if } i \in \mathcal{D} \end{cases},$$

and  $\mathcal{D}$  is the set of cases.

The asymptotic expansion for the survey estimator using the inverse probability weights is modified as  $\mathbf{B}_{2i} = \mathbf{A}_2^{-1} \mathbf{W}_i$ , where  $\mathbf{A}_2 = \lim - (1/n)(\partial \mathbf{U}_2 / \partial \boldsymbol{\beta}_g)$ , and

$$\mathbf{W}_i = \mathbf{U}_{2i} - \sum_{i \in S \cup \mathcal{D}} \frac{\frac{1}{\pi_i} \Delta_i R_i(T_i) I(i \in S \cup \mathcal{D}) \exp(\hat{\boldsymbol{\gamma}}^T \mathbf{X}_{1i} + \boldsymbol{\beta}_g^T \mathbf{X}_{2i})}{n \mathbf{S}^{(0)}(\boldsymbol{\beta}_g, T_i; \hat{\boldsymbol{\gamma}})} \left\{ \mathbf{X}_{2i} - \frac{\mathbf{S}^{(1)}(\boldsymbol{\beta}_g, T_i; \hat{\boldsymbol{\gamma}})}{\mathbf{S}^{(0)}(\boldsymbol{\beta}_g, T_i; \hat{\boldsymbol{\gamma}})} \right\}.$$

The computation of the covariance matrix of  $\hat{\boldsymbol{\beta}}_g$  follows similarly as in Section 2.3.