

Semiparametric Estimation of Treatment Effect in a Pretest-Posttest Study

Selene Leon, Anastasios A. Tsiatis, and Marie Davidian*

Department of Statistics, North Carolina State University, Raleigh, North Carolina 27695-8203, U.S.A.

**email:* davidian@stat.ncsu.edu

SUMMARY. Inference on treatment effects in a pretest-posttest study is a routine objective in medicine, public health, and other fields. A number of approaches have been advocated. We take a semiparametric perspective, making no assumptions about the distributions of baseline and posttest responses. By representing the situation in terms of counterfactual random variables, we exploit recent developments in the literature on missing data and causal inference, to derive the class of all consistent treatment effect estimators, identify the most efficient such estimator, and outline strategies for implementation of estimators that may improve on popular methods. We demonstrate the methods and their properties via simulation and by application to a data set from an HIV clinical trial.

KEY WORDS: Analysis of covariance; Counterfactuals; Influence function; Inverse probability weighting; Semiparametric model; t-test.

1. Introduction

The pretest-posttest trial is ubiquitous in research in medicine, public health, and numerous other fields. In the usual study, subjects are randomized to one of two treatments (e.g., treatment and control), and the response of interest is ascertained for each at baseline (pretreatment) and follow-up. The objective is to evaluate whether treatment affects follow-up response, with baseline responses serving as a basis for comparison. For instance, many HIV clinical trials focus on comparing treatment effects on viral load or CD4 count after a specified period, with baseline observations on these quantities routinely available.

A number of strategies have been advocated for evaluating treatment effect in this setting, including the two-sample t-test comparing follow-up observations from each group, ignoring baseline, the paired t-test comparing differences between follow-up and baseline, and analysis of covariance procedures applied to follow-up responses, with either baseline response only, or baseline and its interaction with treatment, included as a covariates in a linear model; the last is referred to by Yang and Tsiatis (2001) as ANCOVA I and ANCOVA II, respectively. A nonparametric approach in a different spirit was proposed by Quade (1982); we do not consider this here. The two-sample t-test seems predicated on the assumption that baseline and follow-up responses are uncorrelated, so that no precision is to be gained from baseline information, while the others implicitly rely on an apparent assumption of linear dependence between follow-up and baseline. All are often associated with the assumption of normality.

Several authors (e.g., Brogan and Kutner, 1980; Laird, 1983; Crager, 1987; Stanek, 1988; Stein, 1989; Follmann, 1991) have studied these approaches under various assump-

tions, including normality or equality of variances of baseline and follow-up responses. Despite this work and the widespread interest in this problem, there is still no consensus on what approach is preferable under general conditions; in our experience with HIV clinical trials, the paired t-test is often used, because “interest focuses on differences” with no theoretical justification. An attempt to address this issue was presented by Yang and Tsiatis (2001); they studied the large-sample properties of treatment effect estimators based on the above approaches, but under very general conditions, where only the first and second moments of baseline and follow-up responses exist and may differ, and their joint distribution conditional on treatment may be arbitrary. They also considered a “generalized estimating equation” (GEE) approach (e.g., Singer and Andrade, 1997) where baseline and follow-up data are treated as a multivariate response, with arbitrary mean and covariance matrix. Yang and Tsiatis (2001) showed that all these estimators are consistent and asymptotically normal. The GEE estimator is asymptotically equivalent to that from ANCOVA II and most efficient. When the randomization probability is 0.5 or covariance between baseline, and follow-up responses is the same for both treatment and control, then the ANCOVA I estimator is asymptotically equivalent to ANCOVA II/GEE. If baseline and follow-up responses are uncorrelated, the two-sample t-test estimator achieves the same precision as ANCOVA II, but is inefficient otherwise. The paired t-test is equivalent to ANCOVA II only if the difference between follow-up and baseline is uncorrelated with baseline within each treatment.

As the ANCOVA approaches are derived from a supposed linear relationship between baseline and follow-up, some practitioners are reluctant to use them. Asymptotic equivalence

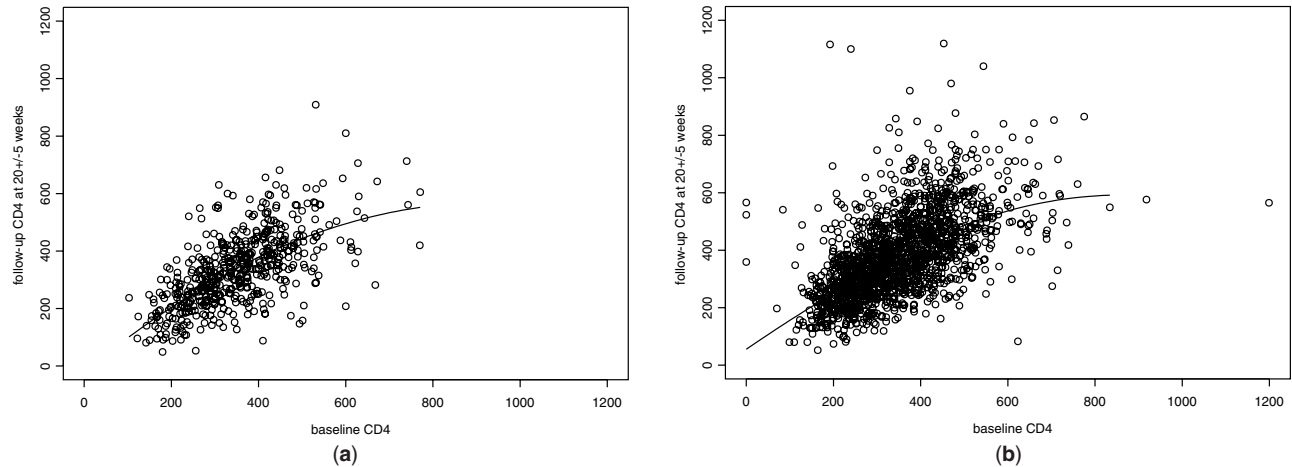


Figure 1. ACTG 175: CD4 counts after 20 ± 5 weeks vs. baseline CD4 counts for patients randomly assigned to (a). ZDV alone (“control”) and (b) the combination of ZDV+ddI, ZDV+ddC, or ddI alone (“treatment”). The solid lines were obtained using the Splus function `loess()` (Cleveland, Grosse, and Shyu, 1993); in (b), use of different span values and deletion of the apparent “high leverage” points with the largest baseline CD4 all lead to similar curvilinear fits.

of the GEE estimator indicates it involves the same considerations. Yang and Tsiatis (2001) showed that consistency and asymptotic normality hold even if linearity is violated. However, as their study restricted attention to such “linear” estimators, it did not address whether or how it might be possible to improve upon these approaches under deviations from linearity, and without limiting distributional assumptions.

Such deviations are commonplace, as illustrated by data from the AIDS Clinical Trials Group (ACTG) protocol: 175 involving 2467 HIV-infected subjects randomized originally to four treatment groups, zidovudine (ZDV) alone, ZDV plus didanosine (ddI), ZDV plus zalcitabine (ddC), or ddI alone in approximately equal numbers (Hammer et al., 1996). Analysis of the primary endpoint of time to progression to AIDS or death showed ZDV to be inferior to the other three therapies, which showed no differences. Figure 1 plots CD4 counts at 20 ± 5 weeks, a follow-up measure that reflects early response to treatment subsequent to the often-observed initial rise in CD4 (e.g., Tsiatis, DeGruttola, and Wulfsohn, 1995), versus baseline CD4 for the ZDV-only (“control”), and other therapies combined (“treatment”) groups; this suggests a possible departure from a straight-line relationship. Indeed, nonlinear relationships are a routine feature of biological phenomena, as is nonnormality. Histograms of CD4 counts at baseline and follow-up for each group (not shown) exhibit the usual asymmetry that motivates the standard analysis on the log, fourth, or cube-root scale.

Data such as those from ACTG 175 highlight the need for inferential methods for the pretest-posttest problem that do not rest on apparent linearity or restrictive distributional assumptions. In this article, we take a new approach, casting the problem in terms of counterfactuals in Section 2, making no assumptions on the joint distribution beyond first moments nor on the relationship between baseline and follow-up response. This suggests an analogy to missing data problems in Section 3, that allows us to exploit the theory of Robins, Rotnitzky, and Zhao (1994), to characterize the class of all consistent treatment effect estimators under

these assumptions and identify the efficient member of the class. In Section 4, we discuss practical implementation of such estimators that should offer improved efficiency over the “popular” methods above. In Section 5, we demonstrate the potential for such improved performance, applying our methods to the ACTG 175 data in Section 6. A fundamental contribution is the definition of a general framework in which the problem may be viewed that provides a strategy for choosing estimators that accommodate features such as nonlinearity.

2. Semiparametric Model Based on Counterfactuals

Let n be the total number of subjects in the trial, each randomized to “control” or “treatment,” with known probabilities $(1 - \delta)$ and δ . Accordingly define $Z_i = 0$ or 1 , for subject i , respectively. Let Y_{1i} and Y_{2i} be i ’s observed baseline and follow-up responses, leading to observed data for i (Y_{1i}, Y_{2i}, Z_i); the subscript i is suppressed when no ambiguity will result.

We develop the model by conceptualizing the situation in terms of counterfactuals or potential outcomes, a key device in the study of causal inference (e.g., Holland, 1986), and then expressing the observable data in terms of these quantities. The variables Y_1 and Z represent phenomena prior to treatment action, while Y_2 is a post-treatment characteristic. Thus, let $Y_2^{(0)}, Y_2^{(1)}$ be the follow-up responses a subject potentially would exhibit if assigned to control and treatment, respectively. The full set of counterfactual random variables $(Y_2^{(0)}, Y_2^{(1)})$ is obviously not observable for any subject, but rather represents what potentially might occur at follow-up under both treatments—including that “counter to the fact” of what might actually be assigned in the trial. We place no restrictions on the joint distribution of the counterfactuals and Y_1 , such as equal variance or independence. We define $\mu_1 = E(Y_1)$, $\sigma_{11} = \text{var}(Y_1)$; for $c = 0, 1$, $\mu_2^{(c)} = E(Y_2^{(c)})$, $\sigma_{22}^{(c)} = \text{var}(Y_2^{(c)})$ and $\sigma_{12}^{(c)} = \text{cov}(Y_1, Y_2^{(c)})$. Thus, e.g., $\mu_2^{(0)} = E(Y_2^{(0)})$ denotes mean follow-up response if all subjects in the population were assigned to control. It is natural to assume that

the observed follow-up response under the subject's actual, assigned treatment corresponds to what potentially would be seen if the subject were assigned to that treatment, i.e., $Y_2 = Y_2^{(0)}(1 - Z) + Y_2^{(1)}Z$. The observed assignment Z is made at random, without regard to baseline status or prognosis. Thus, assume Z is independent of $(Y_1, Y_2^{(0)}, Y_2^{(1)})$. As usual, we assume $(Y_{1i}, Y_{2i}^{(0)}, Y_{2i}^{(1)}, Z_i)$, and hence (Y_{1i}, Y_{2i}, Z_i) are independent and identically distributed (i.i.d.) across i .

Interest focuses on the difference in population mean follow-up response which, under the usual causal inference perspective (Holland, 1986), may be thought of as the difference in means if all subjects in the population were assigned to control or treatment, respectively, i.e., $\beta = \mu_2^{(1)} - \mu_2^{(0)} = E(Y_2^{(1)}) - E(Y_2^{(0)})$. Under our assumptions, in fact, $\beta = E(Y_2|Z=1) - E(Y_2|Z=0)$, the usual expression for the difference of interest in a randomized trial. Also, $E(Y_2|Z) = \mu_2 + \beta Z$ and $E(Y_1|Z) = \mu_1$; we write $\mu_2 = \mu_2^{(0)} = E(Y_2^{(0)})$ for brevity.

The advantage of this framework and expression of β in terms of counterfactual means is that it reveals a structure analogous to that in missing data problems. As we are interested in the difference of the two marginal quantities $\mu_2^{(1)}$ and $\mu_2^{(0)}$, we may view estimation of each mean separately, without reference to the joint distribution of $Y_2^{(0)}, Y_2^{(1)}$; thus, if we identify estimators for $\mu_2^{(1)}$ and $\mu_2^{(0)}$ that are "optimal" in some sense, the "optimal" estimator for β may be obtained as their difference. Accordingly, focus on $\mu_2^{(1)}$; considerations for $\mu_2^{(0)}$ are similar. If we could observe the "full data" $(Y_1, Y_2^{(1)}, Z)$ for all n subjects, then we would estimate $\mu_2^{(1)}$ by the sample mean $n^{-1} \sum_{i=1}^n Y_{2i}^{(1)}$. However, we only observe $Y_2^{(1)}$ for subjects with observed assignment $Z = 1$, so that $Y_2^{(1)}$ is "missing" for subjects with $Z = 0$; i.e., we only observe $(Y_1, ZY_2^{(1)}, Z)$, where $Y_2^{(1)}$ is observed with probability $P(Z=1|Y_1, Y_2^{(1)}) = P(Z=1) = \delta$, so that $Y_2^{(1)}$ is "missing completely at random" (MCAR) (Rubin, 1976). Thus, the two-sample t-test estimator based on observed sample averages,

$$\begin{aligned} n_1^{-1} \sum_{i=1}^n Z_i Y_{2i} - n_0^{-1} \sum_{i=1}^n (1 - Z_i) Y_{2i}, \\ n_1 = \sum_{i=1}^n Z_i, \quad n_0 = \sum_{i=1}^n (1 - Z_i) \end{aligned} \quad (1)$$

may be regarded as a "complete case" estimator for β . As such, while it may be unbiased for β under MCAR, it is likely inefficient, as it doesn't take into account observations on Y_1 . This suggests that a more refined approach to missing data problems may lead to improved estimators that exploit information in Y_1 and the interrelationships among observed variables.

3. A Class of Estimators

3.1 Influence Functions for β

Robins et al. (1994) developed a general large-sample theory of estimation in semiparametric models where data are missing at random (so, including MCAR). They described the class of all (regular, that is, excluding "pathological" cases) estimators under these conditions; they characterized

the form of the influence functions of all members of the class of asymptotically linear estimators (e.g., Newey, 1990) that are consistent and asymptotically normal under general conditions. An estimator $\hat{\beta}$ for β is asymptotically linear with influence function \mathcal{I} if $n^{1/2}(\hat{\beta} - \beta) = n^{-1/2} \sum_{i=1}^n \mathcal{I}_i + o_p(1)$ and $E(\mathcal{I}) = 0, E(\mathcal{I}^T \mathcal{I}) < \infty$, and the asymptotic variance of $\hat{\beta}$ is the variance of the influence function. Robins et al. (1994) identified the most efficient regular, asymptotically linear (RAL) estimator, namely, the one whose influence function has smallest variance. We apply this general theory to our model, to characterize all estimators for β via their influence functions. This allows us not only to demonstrate that the two-sample t-test, paired t-test, and ANCOVA I and II estimators are inefficient members of this class, but also to elucidate the form of the most efficient estimator for β .

If for each $c = 0, 1$, we were able to observe the "full data" $(Y_1, Y_2^{(c)}, Z)$ for all subjects, we would estimate $\mu_2^{(c)}$ by the sample mean $n^{-1} \sum_{i=1}^n Y_{2i}^{(c)}$, with influence function $\varphi^{(c)}(Y_2^{(c)}) = Y_2^{(c)} - \mu_2^{(c)}$. Because $\mu_2^{(c)}, c = 0, 1$, is an explicit function of the distribution of $Y_2^{(c)}$, which we take to be unrestricted, $\varphi^{(0)}$ and $\varphi^{(1)}$ are the only such "full data" influence functions for estimators for $\mu_2^{(0)}$ and $\mu_2^{(1)}$ (Newey, 1990). Of course, we only observe (Y_1, Y_2, Z) for each subject; thus, we require estimators that may be expressed in terms of these quantities. By analogy to missing data problems, it follows from the theory of Robins et al. (1994) that all RAL estimators for β based on the observed data have an influence function of form

$$\begin{aligned} & \left\{ \frac{Z\varphi^{(1)}(Y_2)}{\delta} + \frac{(Z-\delta)h^{(1)}(Y_1)}{\delta} \right\} \\ & - \left[\frac{(1-Z)\varphi^{(0)}(Y_2)}{1-\delta} + \frac{\{(1-Z)-(1-\delta)\}h^{(0)}(Y_1)}{1-\delta} \right], \quad (2) \end{aligned}$$

where $h^{(0)}, h^{(1)}$ are arbitrary functions such that $\text{var}\{h^{(c)}(Y_1)\} < \infty, c = 0, 1$; (2) is the difference of the forms of all observed data influence functions for estimators for $\mu_2^{(1)}$ and $\mu_2^{(0)}$. For arbitrary h with $\text{var}\{h(Y_1)\} < \infty$, we may rewrite (2) as

$$Z(Y_2 - \mu_2 - \beta)/\delta - (1 - Z)(Y_2 - \mu_2)/(1 - \delta) + (Z - \delta)h(Y_1). \quad (3)$$

Thus, (3) characterizes all consistent estimators for β , and we expect that the influence functions for the "popular" estimators in Section 1 may be represented in this form.

3.2 Popular Estimators

The two-sample t-test estimator for β is given in (1). The paired t-test estimator is $\overline{D}_1 - \overline{D}_0$, where $\overline{D}_1 = n_1^{-1} \sum_{i=1}^n Z_i(Y_{2i} - Y_{1i})$ and $\overline{D}_0 = n_0^{-1} \sum_{i=1}^n (1 - Z_i)(Y_{2i} - Y_{1i})$. As in Yang and Tsiatis (2001), the ANCOVA I estimator for β is obtained by ordinary least squares (OLS) regression of Y_2 on (Y_1, Z) , and the ANCOVA II estimator is obtained by OLS regression of $(Y_2 - \overline{Y}_2)$ on $\{(Y_1 - \overline{Y}_1), (Z_i - \overline{Z}), (Y_1 - \overline{Y}_1)(Z_i - \overline{Z})\}^T$. It is straightforward to show that all of these estimators have influence functions of the form

$$\begin{aligned} & Z(Y_2 - \mu_2 - \beta)/\delta - (1 - Z)(Y_2 - \mu_2)/(1 - \delta) \\ & + (Z - \delta)(Y_1 - \mu_1)\eta, \end{aligned} \quad (4)$$

where $\eta = 0, -1/\{\delta(1 - \delta)\}, -\{\delta\sigma_{12}^{(1)} + (1 - \delta)\sigma_{12}^{(0)}\}/\{\sigma_{11}\delta(1 - \delta)\}$, and $-\{(1 - \delta)\sigma_{12}^{(1)} + \delta\sigma_{12}^{(0)}\}/\{\sigma_{11}\delta(1 - \delta)\}$ for the two-sample t-test, paired t-test, ANCOVA I, and ANCOVA II estimators, respectively. Thus, from (4), these estimators are all in class (3), with $h(Y_1) = \eta(Y_1 - \mu_1)$; hence they are consistent and asymptotically normal.

The observations of Yang and Tsiatis (2001) are immediate: if $\delta = 0.5$, then η is identical for ANCOVA I and II and these estimators are asymptotically equivalent. If $\sigma_{12}^{(c)} = 0$, $c = 0, 1$, i.e., uncorrelated baseline and follow-up, then $\eta = 0$ for ANCOVA I/II and both are asymptotically equivalent to the two-sample t-test. Finally, if $Y_2^{(c)} - Y_1$ is uncorrelated with Y_1 , so $\sigma_{12}^{(c)} = \sigma_{11}$, $c = 0, 1$, the paired t-test, is equivalent to ANCOVA I/II. Interestingly, while η values for ANCOVA I/II both involve a “weighted average” of $\sigma_{12}^{(0)}$ and $\sigma_{12}^{(1)}$, the “weighting” for the latter seems counterintuitive, in that $(1 - \delta) = P(Z = 0)$ and $\delta = P(Z = 1)$ are the coefficients of the covariances for treatments 1 and 0, whereas one might expect the reverse.

Not only are all the “popular” estimators in class (3), they belong to the subclass with h a linear function of Y_1 , suggesting there is room for improvement via more general h .

3.3 The Most Efficient Estimator

The most efficient estimator defined by (3) is that attaining the semiparametric efficiency bound; i.e., with an influence function of form (3) with the smallest variance. Identifying this influence function is not straightforward, as it requires finding the infinite-dimensional function h minimizing the variance. The theory of Robins et al. (1994) provides a general procedure for this problem; applying this, the influence function of the most efficient estimator for β is the difference of the most efficient influence functions for $\mu_2 + \beta$ and μ_2 , given by

$$\begin{aligned} & \left[\frac{Z(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(Z - \delta)\{E(Y_2^{(1)}|Y_1) - \mu_2 - \beta\}}{\delta} \right] \\ & - \left[\frac{(1 - Z)(Y_2 - \mu_2)}{1 - \delta} + \frac{(Z - \delta)\{E(Y_2^{(0)}|Y_1) - \mu_2\}}{1 - \delta} \right] \quad (5) \\ & = \frac{Z(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(1 - Z)(Y_2 - \mu_2)}{1 - \delta} - (Z - \delta) \\ & \times \left[\frac{(1 - \delta)\{E(Y_2^{(1)}|Y_1) - \mu_2 - \beta\} + \delta\{E(Y_2^{(0)}|Y_1) - \mu_2\}}{\delta(1 - \delta)} \right], \quad (6) \end{aligned}$$

where, under our assumptions, $E(Y_2^{(1)}|Y_1) = E(Y_2|Y_1, Z = 1)$ and $E(Y_2^{(0)}|Y_1) = E(Y_2|Y_1, Z = 0)$. A sketch of the derivation of (6) is in the Appendix.

For the “popular” estimators, restriction to the subclass with $h(Y_1) = \eta(Y_1 - \mu_1)$ imposes a finite-dimensional structure on h . It is thus straightforward to find the optimal estimator within this subclass by finding the value of η that minimizes the variance of (3) when $h(Y_1) = \eta(Y_1 - \mu_1)$, given

by $\eta = -\{(1 - \delta)\sigma_{12}^{(1)} + \delta\sigma_{12}^{(0)}\}/\{\sigma_{11}\delta(1 - \delta)\}$. Thus, in general, the ANCOVA II estimator for β is most efficient if attention is restricted to estimators having influence function with linear $h(Y_1)$; this is consistent with Yang and Tsiatis (2001).

It is in fact possible to extend these developments to incorporate baseline covariates X_0 to improve efficiency of estimation of β . Considering the “full data” for $c = 0, 1$ to be $(X_0, Y_1, Y_2^{(c)}, Z)$, and applying the Robins et al. (1994) theory, all RAL estimators for β based on the observed data (X_0, Y_1, Y_2, Z) have influence functions of the form (3), with h now a function of (X_0, Y_1) ; the optimal such influence function is of form (6), with $E(Y_2^{(c)}|Y_1)$ replaced by $E(Y_2^{(c)}|X_0, Y_1) = E(Y_2|X_0, Y_1, Z = c)$ for $c = 0, 1$.

4. Practical Implementation

To exploit (6), one must identify estimators for β with this influence function, which will then be “optimal,” as described above. This may be accomplished by finding estimators for $\mu_2^{(1)} = \mu_2 + \beta$ and $\mu_2^{(0)} = \mu_2$ with the influence functions in (5) and taking their difference. An obvious complication is the involvement of the unknown conditional expectations $E(Y_2|Y_1, Z = c)$, $c = 0, 1$, which depend on the unspecified joint distribution of the observed data; thus, a way of deducing these quantities is required. One might use nonparametric smoothing to estimate $E(Y_2|Y_1, Z = c)$, or fit specific parametric models based on inspection of plots like those in Figure 1 (Robins et al., 1994). Nonparametric estimators typically do not attain usual parametric $n^{-1/2}$ -convergence rates, raising concern that effects of smoothing will degrade performance of the estimator for β in small samples, relative to that achieved with a correct, parametric specification. As the results of Newey (1990, pp. 118–119) imply, estimation of $E(Y_2|Y_1, Z = c)$ at a rate faster than $n^{-1/4}$ does not affect the large sample properties of the estimator relative to knowing $E(Y_2|Y_1, Z = c)$; for larger n , such smoothing may be a viable alternative. However, if baseline covariates X_0 are incorporated, multidimensional smoothing is required, which may be prohibitive if $\dim(X_0)$ is large. On the other hand, although the estimator for β will still be consistent, and asymptotically normal as a member of the general class if the chosen parametric form is incorrect, the resulting estimator no longer need have the optimal influence function, so could in fact be inferior to “popular” estimators. Choosing a parametric model can be tricky; for the ACTG 175 data, the nature of the “true” relationship is ambiguous. Figure 1 suggests several plausible parametric models for each group, e.g., a linear, quadratic, or nonlinear (exponential) function. Regardless, one is still faced with the issue of deriving appropriate estimators for $\mu_2^{(1)}$ and $\mu_2^{(0)}$.

We propose a strategy that may be regarded as a “compromise” between fully nonparametric smoothing and parametric modeling that leads straightforwardly to a general form of an estimator for β that not only will improve on the “popular” ones, but that has the optimal influence function under conditions we elucidate shortly. The approach is based on restricting the search for estimators for β to those with

influence functions of the form

$$\frac{Z_i(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(1 - Z_i)(Y_2 - \mu_2)}{1 - \delta} + (Z - \delta)f^T(Y_1)\alpha, \quad \alpha \in \mathbb{R}^k, \quad (7)$$

where $f(Y_1) = \{f_1(Y_1), \dots, f_k(Y_1)\}^T$ is a k -vector of basis functions. For example, the $(k - 1)$ -order polynomial basis takes $f(Y_1) = \{1, Y_1, Y_1^2, \dots, Y_1^{k-1}\}^T$; alternatively, one may choose a spline basis or discretization basis with $f(Y_1) = \{I(Y_1 < t_1), I(t_1 \leq Y_1 < t_2), \dots, I(t_{k-2} \leq Y_1 < t_{k-1}), I(Y_1 \geq t_{k-1})\}^T$ for $t_1 < t_2 < \dots < t_{k-1}$. Thus, (7) may be viewed as restricting the search for h in (6) to the linear space spanned by $f(Y_1)$. If the basis is sufficiently rich that this space is a good approximation to the space of all possible h , then the resulting estimators should be close to “optimal” if they are “optimal” within the restricted class.

We thus find the most efficient influence function in class (7), which follows by identifying α minimizing the variance of (7). Under our assumptions, the first two terms of (7) are uncorrelated, so this is equivalent to minimizing $\text{var}(A - B^T\alpha)$, where A corresponds to the first two terms and $B = -(Z - \delta)f(Y_1)$. This is an unweighted least squares problem, so that $\alpha^T = \text{cov}(A, B)\{\text{var}(B^T)\}^{-1}$, which may be shown to be $\alpha^T = -\{(1 - \delta)\Sigma_{fY_2}^{(1)T} + \delta\Sigma_{fY_2}^{(0)T}\}\Sigma_{ff}^{-1}/\{\delta(1 - \delta)\}$, where $\Sigma_{fY_2}^{(1)} = E\{f(Y_1)(Y_2 - \mu_2 - \beta)|Z = 1\}$, $\Sigma_{fY_2}^{(0)} = E\{f(Y_1)(Y_2 - \mu_2)|Z = 0\}$, and $\Sigma_{ff} = E\{f(Y_1)f^T(Y_1)\}$. In fact, this α may be found by the sum of separate regressions of each term on B . Thus, the optimal influence function in the class (7) is

$$\begin{aligned} & \left\{ \frac{Z_i(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(Z - \delta)\Sigma_{fY_2}^{(1)T}\Sigma_{ff}^{-1}f(Y_1)}{\delta} \right\} \\ & - \left\{ \frac{(1 - Z)(Y_2 - \mu_2)}{1 - \delta} + \frac{(Z - \delta)\Sigma_{fY_2}^{(0)T}\Sigma_{ff}^{-1}f(Y_1)}{1 - \delta} \right\} \quad (8) \\ & = \frac{Z_i(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(1 - Z)(Y_2 - \mu_2)}{1 - \delta} \\ & - (Z - \delta) \left\{ \frac{(1 - \delta)\Sigma_{fY_2}^{(1)} + \delta\Sigma_{fY_2}^{(0)}}{\delta(1 - \delta)} \right\}^T \Sigma_{ff}^{-1}f(Y_1) \quad (9) \end{aligned}$$

Derivation of an estimator for β with influence function (9) is straightforward: consider the equivalent representation (8) to deduce estimators for each of $\mu_2^{(1)} = \mu_2 + \beta$ and $\mu_2^{(0)} = \mu_2$ and take their difference. An estimator for $\mu_2^{(1)}$ with influence function equal to the first term in braces in (8) may be found by equating the sample average of such terms to zero to obtain $\mu_2^{(1)} = [n^{-1} \sum_{i=1}^n \{Z_i Y_{2i} - (Z_i - \delta)\Sigma_{fY_2}^{(1)T}\Sigma_{ff}^{-1}f(Y_{1i})\}] / \{n^{-1} \sum_{i=1}^n Z_i\}$. This suggests the estimator for $\mu_2^{(1)}$ found by substituting sample moment analogs for each quantity in this expression. Using similar calculations for the second term in braces in (8) to isolate $\mu_2^{(0)}$ and defining $S_{fY_2}^{(c)} = \sum_{i=1}^n I(Z_i = c)f(Y_{1i})(Y_{2i} - \bar{Y}_2^{(c)})$, $c = 0, 1$, $S_{ff} = \sum_{i=1}^n f(Y_{1i})f^T(Y_{1i})$, and $S_{fZ} = \sum_{i=1}^n (Z_i - n_1/n)f(Y_{1i})$, by taking the difference, we

obtain the estimator for β

$$\hat{\beta} = \bar{Y}_2^{(1)} - \bar{Y}_2^{(0)} - n \left(\frac{S_{fY_2}^{(1)}}{n_1^2} + \frac{S_{fY_2}^{(0)}}{n_0^2} \right)^T S_{ff}^{-1} S_{fZ}, \quad (10)$$

which may be shown to have influence function (9). It is straightforward to show that the variance of (9), and hence the large-sample variance of $n^{1/2}(\hat{\beta} - \beta)$, is given by

$$\frac{\sigma_{22}^{(1)}}{\delta} + \frac{\sigma_{22}^{(0)}}{1 - \delta} - \delta(1 - \delta) \left(\frac{\Sigma_{fY_2}^{(1)}}{\delta} + \frac{\Sigma_{fY_2}^{(0)}}{1 - \delta} \right)^T \Sigma_{ff}^{-1} \left(\frac{\Sigma_{fY_2}^{(1)}}{\delta} + \frac{\Sigma_{fY_2}^{(0)}}{1 - \delta} \right). \quad (11)$$

This suggests the estimator for sampling variance of $\hat{\beta}$, given by

$$\frac{S_{22}^{(1)}}{n_1^2} + \frac{S_{22}^{(0)}}{n_0^2} - n_1 n_0 \left(\frac{S_{fY_2}^{(1)}}{n_1^2} + \frac{S_{fY_2}^{(0)}}{n_0^2} \right)^T S_{ff}^{-1} \left(\frac{S_{fY_2}^{(1)}}{n_1^2} + \frac{S_{fY_2}^{(0)}}{n_0^2} \right), \quad (12)$$

where $S_{22}^{(c)} = \sum_{i=1}^n I(Z_i = c)(Y_{2i} - \bar{Y}_2^{(c)})^2$, $c = 0, 1$.

A modification is to use different sets of basis functions, $f^{(c)}(Y_1)$, $c = 0, 1$, say, for each term in (8). Alternatively, these developments suggest applying a similar approach to (5) and (6), representing $E(Y_2|Y_1, Z = c)$, $c = 0, 1$, by linear combinations of basis functions. It is straightforward to show that this leads to the same class of estimators represented by (7), suggesting that, in practice, insight may be gained into the choice of basis by examining plots such as Figure 1. Thus, the approach may be viewed as intermediate to completely nonparametric and fully parametric estimation of the $E(Y_2|Y_1, Z = c)$, approximating these quantities by a finite-dimensional, flexible form emphasizing the predominant trend apparent in the data. From the derivation of the optimal α , to achieve efficiency gains, it is essential that this be carried out by unweighted regression, even if $\text{var}(Y_2|Y_1, Z = c)$, $c = 0, 1$, is not constant with respect to Y_1 . When the true conditional expectations follow such a form, and k and the chosen basis functions represent them exactly correctly, the method yields asymptotically the most efficient estimator for β ; otherwise, the estimator is still consistent and asymptotically normal, and we expect gains over the “popular” estimators as long as the basis approximates a broad range of relationships. That is, the estimator is locally semiparametric efficient (e.g., Robins et al., 1994, Section 4.1). Moreover, for a given choice of basis, (12) consistently estimates the true variance of the estimator, regardless of whether $E(Y_2|Y_1, Z = c)$ can be represented by a linear combination of the chosen basis functions. In the next section, we demonstrate that close-to-“optimal” inference is obtained when the true relationship is nonlinear but the basis representation captures its salient features.

The influence function and its variance depend on moments of squares and crossproducts of elements of $f(Y_1)$ through Σ_{ff} , and the covariances of $Y_2^{(c)}$, $c = 0, 1$, with elements of $f(Y_1)$ through $\Sigma_{fY_2}^{(c)}$, which are estimated in (10) and (12) by sample analogs. Thus, e.g., the quadratic basis $f(Y_1) = (1, Y_1, Y_1^2)^T$ used in the next section involves estimation

of not only σ_{11} , $\sigma_{12}^{(c)}$, and $\sigma_{22}^{(c)}$, but also of $\text{cov}(Y_1^2, Y_2^{(c)})$, $c = 0, 1$, and the coefficients of skewness and excess kurtosis of the distribution of Y_1 . This suggests that attempting to gain efficiency over “popular” estimators in this way with small sample sizes is unwise. However, in situations such as large clinical trials, this approach may be fruitful. The simulation evidence in Section 5 indicates that impressive efficiency gains are possible in the moderate-to-large sample sizes where the improved estimators are expected to perform well.

When baseline covariates are available, the strategy is immediately extended by replacing $f(Y_1)$ in (7) by a k -vector of basis functions $f(X_0, Y_1)$; e.g., one might choose a polynomial basis including interactions of powers of Y_1 with elements of X_0 .

5. Simulation Evidence

We carried out several simulation studies, and report here on results for four scenarios, each involving 5000 Monte Carlo replications, $\beta = 0.5$, $\delta = 0.5$, and $Y_1 \sim \mathcal{N}(\mu_1 = 0, \sigma_{11} = 1)$. We estimated β by 1) ANCOVA I and II, 2) the two-sample t-test estimator, 3) the paired t-test estimator, 4) (10), with quadratic polynomial basis, denoted QUAD, and 5) the estimator formed by estimating $E(Y_2|Y_1, Z = c)$, $c = 0, 1$, via locally weighted polynomial smoothing using `proc loess` in SAS (SAS Institute, 2000) with quadratic polynomials, substituting in (5), finding estimators for $\mu_2^{(0)}$ and $\mu_2^{(1)}$ by equating sample averages of each term in (5) to zero, and taking their difference, denoted LOESS. For “popular” estimators, standard errors were obtained, both by substituting sample moments in the asymptotic variance formulae suggested by their influence functions (given explicitly in Section 2 of Yang and Tsiatis, 2001), and using the “usual” expressions for these estimators; e.g., for ANCOVA I and II obtained from standard OLS formulae. For QUAD, standard errors were obtained from (12); for LOESS, standard errors were computed by taking the variance of (6), assuming $E(Y_2|Y_1, Z = c)$ known and substituting sample moment quantities. Nominal 95% Wald confidence intervals for β were constructed as the estimate ± 1.96 times the asymptotic-formula standard error.

Follow-up responses for the first two scenarios were generated from the quadratic model

$$Y_{2i} = (\mu_2 + \beta Z_i) + \beta_1(Y_{1i} - \mu_1) + \beta_2\{(Y_{1i} - \mu_1)^2 - \sigma_{11}\} + \epsilon_i, \\ \epsilon_i \sim \mathcal{N}((0, 1)) \quad (13)$$

with $\mu_2 = -0.25$. Normality of baseline and follow-up may be considered the “most favorable” distribution for the “popular” estimators, which are often thought to be predicated on normality. However, from Section 3.2, these estimators are consistent more generally. Situation Q1 is based on (13) with $(\beta_1, \beta_2) = (0.5, 0.4)$, yielding a discernible curvilinear relationship between baseline and follow-up, as depicted in Figure 2(a). The “popular” estimators assume a linear relationship; thus, the linear correlation of 0.40 between baseline and follow-up in each group is of interest. The second situation, Q2, exemplified in Figure 2(b), with $(\beta_1, \beta_2) = (0.1, 0.1)$, involves a “low” correlation of 0.10 in each group. Table 1 shows results for $n = 100, 500$, which for $n = 100$, are similar

to those of Yang and Tsiatis (2001); they used (13), but with interactions between baseline and treatment; they also mislabelled Monte Carlo standard deviation and standard error estimates as “variance.” In all cases, bias is negligible, standard error estimates are reliable, and coverage probabilities are close to the nominal level. ANCOVA I/II are equivalent; for Q2, the two-sample t-test performs well, and the paired t-test is inefficient for both scenarios, as expected. Most striking is the considerable gain in efficiency attained by QUAD over the “popular” approaches in the realistic situation of Q1, with a moderate degree of curvilinear association. The LOESS method is virtually equivalent to QUAD for $n = 500$, in both Q1 and Q2. However, for the smaller $n = 100$, it exhibits some loss of efficiency, perhaps reflecting the concern raised in Section 4. Overall, the proposed approach, which seeks to enhance performance by exploiting the nature of the relationship between baseline and follow-up, can dramatically improve precision. From Table 1, as expected, little is to be gained when this relationship is weak (Q2).

In Q1 and Q2, the basis functions coincided with the true form of $E(Y_2|Y_1, Z)$. To investigate performance for a more complicated relationship, we generated data from

$$Y_{2i} = \beta_0 + \beta Z_i + e^{\beta_1 + \beta_2 Y_{1i}} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, 1), \quad (14)$$

where now μ_2 depends on $(\beta_0, \beta_1, \beta_2)$. In the first situation, N1, $(\beta_0, \beta_1, \beta_2) = (-4.0, 1.0, 0.5)$. This results in curvature typified by Figure 2(c), with “high” linear correlation between baseline and follow-up of about 0.80 in each group. Situation N2, with $(\beta_0, \beta_1, \beta_2) = (-4.0, 1.4, 0.1)$, produces haphazard scatterplots, as in Figure 2(d), and “weak” correlation of roughly 0.30. In (14), $E(Y_2|Y_1, Z)$ is not quadratic; thus, as an “ideal” benchmark, we also estimated β by taking the true $E(Y_2|Y_1, Z)$ to be known and finding the optimal estimator based on (6); this is denoted as BENCHMARK in Table 2. Standard errors for this estimator were calculated in a manner similar to those for LOESS. For N1, LOESS and QUAD perform well relative to the unachievable “ideal,” and offer appreciable gains in efficiency relative to the “popular” estimators. This is despite “high” correlation, although again LOESS is less precise when $n = 100$. Not unexpectedly, for N2, with no discernible relationship, the BENCHMARK, QUAD, and LOESS estimators offer no improvement over the best “linear” estimators.

It is natural to wonder if these efficiency gains translate to increased power for testing the usual hypothesis $H_0: \beta = 0$. Table 3 shows empirical size and power for Wald tests of H_0 under scenario N1. The tests achieve the nominal level, most exhibiting some elevation when $n = 100$. The proposed approach yields 10–25% increases in power over the nearest competitors, except for $n = 500$ with large alternatives.

6. Treatment Effect in ACTG 175

Figure 1 shows a possibly nonlinear baseline-follow-up relationship in each ACTG 175 group (with correlations of 0.55 and 0.65 for control and treatment); here, $\delta = 0.75$. It is standard to analyze CD4 counts on a scale where they appear symmetrically distributed, to achieve the approximate normality widely thought required for valid inference via

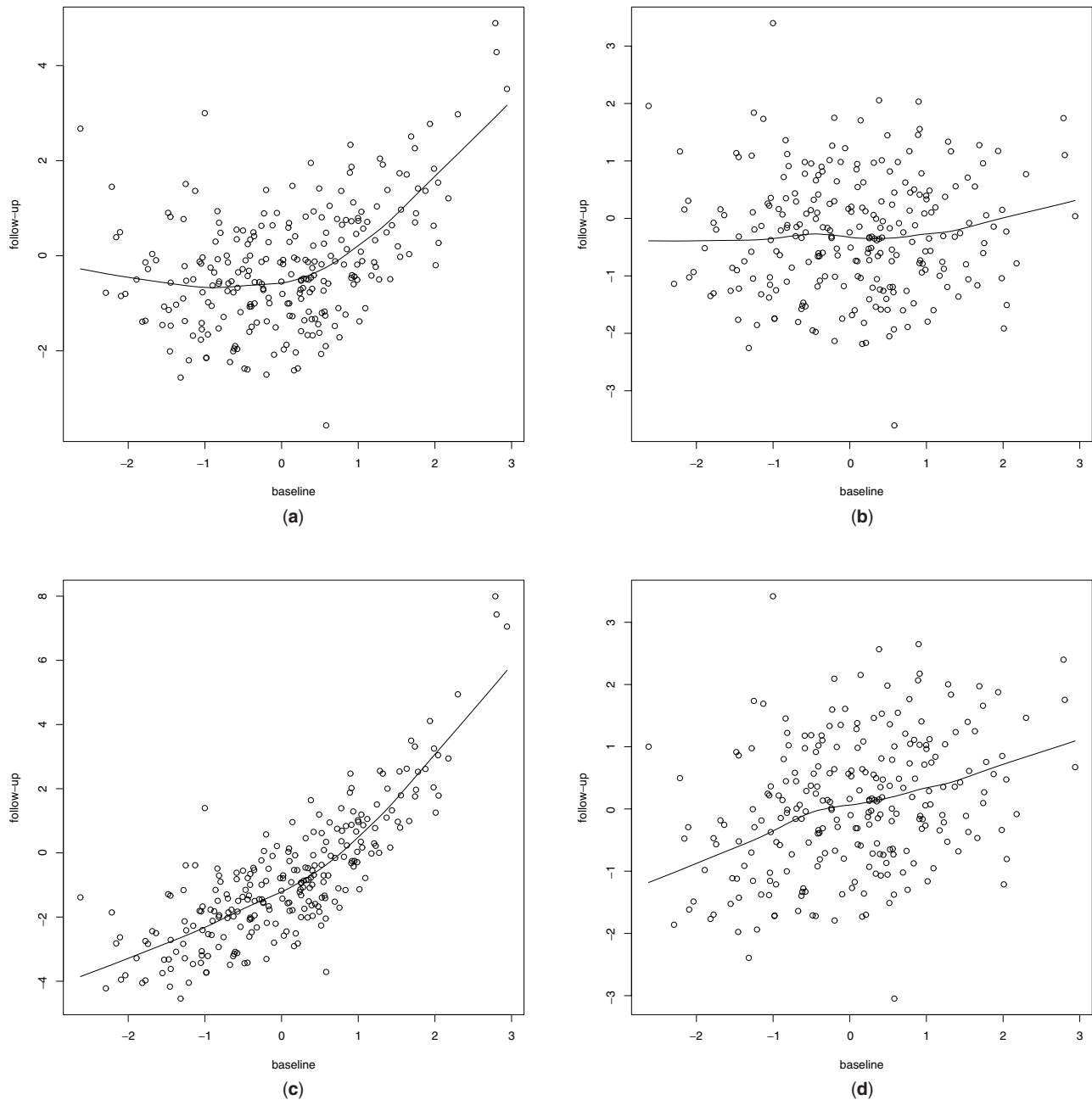


Figure 2. Simulated data for $n = 500$ from scenarios (a) Q1 and (b) Q2, which are both based on (13), and scenarios (c) N1 and (d) N2, which are both based on (14). Smooth fits (solid line) were obtained using the Splus function `loess()` (Cleveland et al., 1993). Each panel depicts data for roughly 250 subjects randomized to control.

“popular” methods. However, Section 3 shows this is not needed for consistency. If interest focuses on mean CD4, inference under a monotone transformation to normality instead addresses median CD4.

Table 4 presents results using the approaches studied in Section 5, for data from 2130 subjects, excluding those missing baseline or 20 ± 5 week CD4. Also given are two estimators incorporating baseline covariates $X_0 = \{\text{weight (kg), Karnofsky score (0-100), days of prestudy antiretroviral ther-$

apy, symptomatic status (0/1), IV drug use (0/1)}—one using basis functions $f(X_0, Y_1) = (1, Y_1, Y_1^2, X_{01}, \dots, X_{05})^T$ and the other obtained by fitting $E(Y_2|X_0, Y_1, Z = c)$, $c = 0, 1$, via generalized additive models (Hastie and Tibshirani, 1990). The proposed methods have the smallest standard errors, and incorporation of baseline covariates yields further improvement. The gains are slight, however, likely owing to the weak curvature in Figure 1 and weak prognostic effect of the covariates.

Table 1

Simulation results for true quadratic relationship (13), 5000 Monte Carlo data sets. Estimators and scenarios are as described in the text.

Estimator	MC mean	MC SD	Asymp. SE	OLS SE	MSE ratio	CP	<i>n</i>
Scenario Q1, “moderate” association							
LOESS	0.504	0.218	0.202	–	0.89	0.93	100
QUAD	0.507	0.205	0.200	–	1.00	0.93	100
ANCOVA II	0.506	0.237	0.228	0.231	0.75	0.94	100
ANCOVA I	0.506	0.233	0.228	0.231	0.77	0.94	100
Paired t-test	0.508	0.255	0.251	0.251	0.65	0.95	100
Two sample t-test	0.506	0.254	0.251	0.251	0.65	0.94	100
LOESS	0.501	0.090	0.091	–	0.98	0.96	500
QUAD	0.501	0.089	0.089	–	1.00	0.95	500
ANCOVA II	0.502	0.103	0.103	0.103	0.75	0.95	500
ANCOVA I	0.502	0.102	0.103	0.103	0.76	0.95	500
Paired t-test	0.501	0.111	0.112	0.112	0.64	0.95	500
Two sample t-test	0.503	0.112	0.112	0.112	0.63	0.95	500
Scenario Q2, “weak” association							
LOESS	0.506	0.214	0.189	–	0.90	0.91	100
QUAD	0.505	0.203	0.199	–	1.00	0.94	100
ANCOVA II	0.505	0.205	0.202	0.205	0.98	0.94	100
ANCOVA I	0.505	0.204	0.202	0.204	0.99	0.94	100
Paired t-test	0.511	0.269	0.272	0.272	0.57	0.95	100
Two sample t-test	0.505	0.204	0.204	0.204	0.99	0.95	100
LOESS	0.499	0.091	0.089	–	0.97	0.94	500
QUAD	0.499	0.089	0.089	–	1.00	0.95	500
ANCOVA II	0.499	0.091	0.090	0.091	0.98	0.95	500
ANCOVA I	0.499	0.090	0.090	0.090	0.98	0.95	500
Paired t-test	0.499	0.121	0.121	0.121	0.55	0.95	500
Two sample t-test	0.499	0.091	0.091	0.091	0.97	0.95	500

MC Mean is Monte Carlo average, MC SD is Monte Carlo standard deviation, Asymp. SE is the average of estimated standard errors based on the asymptotic theory (see text), OLS SE is the average of estimated standard errors based on OLS for the “popular” estimators, MSE Ratio is Mean Square Error (MSE) for QUAD divided by MSE of the indicated estimator, CP is empirical coverage probability of confidence interval using asymptotic SEs, and *n* is the total sample size.

7. Discussion

Casting the pretest-posttest problem in a counterfactual framework, we have exploited advances in the missing data literature to characterize a general class of estimators for treatment effect. We have shown that “popular” approaches are inefficient, and that via our approach, improvement is possible by taking into account the nature of the relationship between baseline and follow-up response. We do not recommend the proposed estimators with very small samples, as they require estimation of features that may be not be well-identified under these conditions. The formulation also yields an approach for incorporating baseline covariate information to further improve efficiency. Our strategy differs from that of directly modeling parametrically the relationship between follow-up and baseline and other factors. In this approach, consistency is predicated on correct specification of the model (and possible distributional assumptions), while ours yields consistent estimators of treatment effect regardless of if, e.g., the chosen basis does not accurately reflect the true relationship.

When there exists an interaction between baseline response and treatment, a philosophical issue is whether interest should

indeed focus on main effects. Proponents of the “large, simple trial” (e.g., Friedman, Furberg, and DeMets, 1996, p. 56) downplay interactions and focus on overall effect, while researchers in other settings have a different view. We do not take a position in this debate. The proposed estimators are for such a main effect; whether there is an interaction or not, the methods estimate this effect consistently and exploit relationships among variables solely to gain efficiency.

It is in fact straightforward to extend the development to more than two treatments, and to observational studies where treatment assignment is not random and it is assumed that treatment assignment is independent of prognosis (i.e., follow-up counterfactuals), given baseline covariates. Also, we have assumed that there are no missing data; we purposely have restricted attention to this setting to elucidate the general structure of the problem. Missing data, particularly at follow-up, are a feature of many pretest-posttest studies. The theory of Robins et al. (1994) may be applied to derive the class of all estimators for β when responses are missing at random, to deduce a strategy similar to that here to identify estimators under these conditions. We report on this development elsewhere.

Table 2

Simulation results for true nonlinear relationship (14), 5000 Monte Carlo data sets. Estimators and scenarios are as described in the text.

Estimator	MC mean	MC SD	Asymp. SE	OLS SE	MSE ratio	CP	<i>n</i>
Scenario N1, “high” association							
BENCHMARK	0.507	0.203	0.205	–	1.05	0.94	100
LOESS	0.504	0.218	0.237	–	0.91	0.96	100
QUAD	0.507	0.207	0.201	–	1.00	0.93	100
ANCOVA II	0.506	0.236	0.228	0.231	0.77	0.93	100
ANCOVA I	0.506	0.232	0.228	0.231	0.80	0.94	100
Paired t-test	0.505	0.257	0.254	0.254	0.65	0.94	100
Two sample t-test	0.503	0.387	0.384	0.384	0.29	0.94	100
BENCHMARK	0.501	0.089	0.092	–	1.03	0.96	500
LOESS	0.501	0.090	0.096	–	1.00	0.96	500
QUAD	0.501	0.090	0.090	–	1.00	0.95	500
ANCOVA II	0.502	0.103	0.103	0.103	0.77	0.95	500
ANCOVA I	0.502	0.103	0.103	0.103	0.77	0.95	500
Paired t-test	0.502	0.114	0.114	0.114	0.63	0.95	500
Two sample t-test	0.504	0.173	0.172	0.172	0.27	0.95	500
Scenario N2, “mild” association							
BENCHMARK	0.505	0.201	0.204	–	1.03	0.95	100
LOESS	0.506	0.214	0.191	–	0.90	0.92	100
QUAD	0.505	0.203	0.199	–	1.00	0.94	100
ANCOVA II	0.505	0.203	0.200	0.203	1.00	0.94	100
ANCOVA I	0.505	0.202	0.200	0.202	1.01	0.94	100
Paired t-test	0.509	0.231	0.233	0.233	0.77	0.95	100
Two sample t-test	0.503	0.219	0.217	0.217	0.86	0.95	100
BENCHMARK	0.499	0.089	0.091	–	1.00	0.96	500
LOESS	0.499	0.091	0.089	–	0.97	0.95	500
QUAD	0.499	0.089	0.089	–	1.00	0.95	500
ANCOVA II	0.499	0.090	0.089	0.090	1.00	0.95	500
ANCOVA I	0.499	0.089	0.089	0.090	1.00	0.95	500
Paired t-test	0.499	0.103	0.104	0.104	0.75	0.95	500
Two sample t-test	0.499	0.098	0.097	0.097	0.84	0.95	500

Table 3

Empirical size and power of Wald tests (estimate/asymptotic standard error estimate) of $H_0 : \beta = 0$ under scenario N1 with (14), each based on 5000 Monte Carlo data sets. Empirical size was found by simulations with $\beta = 0$. Empirical power is under the indicated alternative.

<i>n</i>	Size		Power					
	$\beta = 0$		$\beta = 0.25$		$\beta = 0.40$		$\beta = 0.50$	
	100	500	100	500	100	500	100	500
BENCHMARK	0.07	0.05	0.26	0.80	0.52	0.99	0.70	1.00
LOESS	0.04	0.04	0.18	0.76	0.39	0.99	0.56	1.00
QUAD	0.07	0.05	0.25	0.79	0.53	0.99	0.71	1.00
ANCOVA II	0.07	0.05	0.21	0.69	0.43	0.97	0.60	1.00
ANCOVA I	0.06	0.05	0.21	0.69	0.43	0.97	0.60	1.00
Paired t-test	0.06	0.05	0.17	0.61	0.37	0.94	0.52	0.99
Two sample t-test	0.06	0.05	0.11	0.32	0.19	0.66	0.26	0.83

Table 4

Treatment effect estimates for 20 ± 5 week CD4 counts for ACTG 175.

Estimator	$\hat{\beta}$	Asymptotic SE	OLS SE
BASE-GAM	50.001	5.080	–
BASE-QUAD	50.837	4.957	–
LOESS	49.799	5.222	–
QUAD	49.792	5.330	–
ANCOVA II	49.404	5.384	5.842
ANCOVA I	49.313	5.385	5.840
Paired t-test	50.148	5.686	6.109
Two sample t-test	45.506	6.767	7.203

The methods are as denoted as in Tables 1–3; in addition, BASE-QUAD denotes the proposed basis function method, including up to quadratic terms in baseline CD4, and linear terms involving baseline covariates, and BASE-GAM denotes the proposed method estimating the conditional expectations using generalized additive models. Asymptotic SE is estimated standard error based on the influence function, and OLS SE is estimated standard error based on the “usual” approaches for the “popular” estimators as described in Section 5.

ACKNOWLEDGEMENTS

This work was supported by NIH grants AI31789, CA085848, and CA51962. The authors are grateful to Michael Hughes, Heather Gorski, and the ACTG for providing the ACTG 175 data, and to the reviewers for very helpful comments on the original version.

RÉSUMÉ

En médecine, en santé publique et dans d'autres domaines, les études avant-après, sont classiques. Pour réaliser des inférences sur l'effet d'un traitement en ajustant sur la valeur de la variable réponse avant traitement, nous adoptons une approche semiparamétrique qui n'exige aucune hypothèse sur les distributions des réponses aux deux temps de mesure. Nous utilisons des développements récents concernant l'analyse des observations incomplètes et l'inférence causale, et nous utilisons la notion de variable aléatoire contrefactuelle pour dériver la classe de tous les estimateurs convergents de l'effet traitement, identifier le plus efficace et mettre en évidence les stratégies de mise en œuvre de ces estimateurs susceptibles de conduire à des performances meilleures que les méthodes courantes. Nous évaluons les propriétés de ces méthodes ces méthodes à l'aide de simulations, et nous les illustrons à partir des données d'un essai thérapeutique comparant des traitements du VIH.

REFERENCES

- Brogan, D. R. and Kutner, M. H. (1980). Comparative analysis of pretest-posttest research designs. *American Statistician* **34**, 229–232.
- Cleveland, W. S., Grosse, E., and Shyu, W. M. (1993). Local regression models. In *Statistical Models in S*, J. M. Chambers and T. J. Hastie (eds), 309–376. New York: Chapman and Hall.
- Crager, M. R. (1987). Analysis of covariance in parallel-group clinical trials with pretreatment baseline. *Biometrics* **43**, 895–901.
- Follmann, D. A. (1991). The effect of screening on some pretest-posttest test variances. *Biometrics* **47**, 763–771.
- Friedman, L. M., Furberg, C. D., and DeMets, D. L. (1996). *Fundamentals of Clinical Trials*, 3rd edition. St. Louis: Mosby.
- Hammer, S. M., Katzstein, D. A., Hughes, M. D., et al. for the AIDS Clinical Trials Group Study 175 Study Team. (1996). A trial comparing nucleoside monotherapy with combination therapy in HIV-infected adults with CD4 cell counts from 200 to 500 per cubic millimeter. *New England Journal of Medicine* **335**, 1081–1089.
- Hastie, T. J. and Tibshirani, R. J. (1990). *Generalized Additive Models*. London: Chapman and Hall.
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association* **81**, 945–960.
- Laird, N. (1983). Further comparative analyses of pretest-posttest research designs. *American Statistician* **37**, 329–340.
- Newey, W. K. (1990). Semiparametric efficiency bounds. *Journal of Applied Econometrics* **5**, 99–135.
- Quade, D. (1982). Nonparametric analysis covariance by matching. *Biometrics* **38**, 597–611.
- Robins, J. M., Rotnitzky, A., and Zhao, L. P. (1994). Estima-

tion of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association* **89**, 846–866.

- Rubin, D. B. (1976). Inference and missing data. *Biometrika* **63**, 581–592.
- SAS Institute (2000). *SAS/STAT User's Guide*, Version 8, 4th edition. Cary, North Carolina: SAS Publishing.
- Singer, J. M. and Andrade, D. F. (1997). Regression models for the analysis of pretest/posttest data. *Biometrics* **53**, 729–735.
- Stanek, E. J., III (1988). Choosing a pretest-posttest analysis. *American Statistician* **42**, 178–183.
- Stein, R. A. (1989). Adjusting treatment effects for baseline and other predictor variables. In *Proceedings of the ASA Biopharmaceutical Section*, 274–280.
- Tsiatis, A. A., DeGruttola, V., and Wulfsohn, M. S. (1995). Modeling the relationship of survival to longitudinal data measured with error: Applications to survival and CD4 counts in patients with AIDS. *Journal of the American Statistical Association* **90**, 27–37.
- Yang, L. and Tsiatis, A. A. (2001). Efficiency study of estimators for a treatment effect in a pretest-posttest trial. *American Statistician* **55**, 314–321.

Received November 2002. Revised June 2003.

Accepted June 2003.

APPENDIX

Sketch of Derivation of (6)

Following the theory of Robins et al. (1994), the set of all functions of (Y_2, Y_1, Z) of the form $(Z - \delta)h(Y_1)$, $E\{h^2(Y_1)\} < \infty$, is a linear subspace of the Hilbert space of all mean-zero functions $\varphi(Y_1, Y_2, Z)$, with $E\{\varphi^2(Y_1, Y_2, Z)\} < \infty$ and covariance inner product. Denoting this subspace as Λ_2 , the most efficient estimator for β is that with influence function

$$\left\{ \frac{Z(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(1 - Z)(Y_2 - \mu_2)}{(1 - \delta)} \right\} - \Pi \left\{ \frac{Z(Y_2 - \mu_2 - \beta)}{\delta} - \frac{(1 - Z)(Y_2 - \mu_2)}{(1 - \delta)} \middle| \Lambda_2 \right\}, \quad (\text{A.1})$$

where $\Pi(\cdot|\Lambda_2)$ is the projection of the argument onto Λ_2 . Because projection is a linear operation, the projection may be found as the difference of the projections of the components of the first term in (A.1). To find $\Pi\{Z(Y_2 - \mu_2 - \beta)/\delta|\Lambda_2\}$, we must find $h^{(1)}(Y_1)$ such that $E\{[Z(Y_2 - \mu_2 - \beta)/\delta - (Z - \delta)h^{(1)}(Y_1)](Z - \delta)h(Y_1)\} = 0$ for all $h(Y_1)$; i.e., we require that $E\{[Z(Y_2 - \mu_2 - \beta)/\delta - (Z - \delta)h^{(1)}(Y_1)](Z - \delta)Y_1\} = 0$ a.s. This may be written equivalently as $E\{Z(Y_2 - \mu_2 - \beta)(Z - \delta)/\delta|Y_1\} = h^{(1)}(Y_1)E\{(Z - \delta)^2|Y_1\}$ a.s. Using independence of Z and Y_1 , the left-hand side of this expression equals $\{E(Y_2^{(1)}|Y_1) - \mu_2 - \beta\}(1 - \delta)$ and $E\{(Z - \delta)^2|Y_1\} = \delta(1 - \delta)$, so that $h^{(1)}(Y_1) = \{E(Y_2^{(1)}|Y_1) - \mu_2 - \beta\}/\delta$, which yields $\Pi\{Z(Y_2 - \mu_2 - \beta)/\delta|\Lambda_2\} = (Z - \delta)\{E(Y_2^{(1)}|Y_1) - \mu_2 - \beta\}/\delta$. Similarly, $\Pi\{Z(Y_2 - \mu_2)/(1 - \delta)|\Lambda_2\} = (Z - \delta)\{E(Y_2^{(0)}|Y_1) - \mu_2\}/(1 - \delta)$. Substituting into (A.1) yields the result.