# *University of California, Berkeley*
## U.C. Berkeley Division of Biostatistics Working Paper Series

# Targeted Maximum Likelihood Estimation of Conditional Relative Risk in a Semi-parametric Regression Model

Cathy Tuglus[*]      Kristin E. Porter[†]

Mark J. van der Laan[‡]

[*]University of California, Berkeley, ctuglus@gmail.com

[†]University of California, Berkeley, kristinporter@berkeley.edu

[‡]University of California - Berkeley, laan@berkeley.edu

# Targeted Maximum Likelihood Estimation of Conditional Relative Risk in a Semi-parametric Regression Model

Cathy Tuglus, Kristin E. Porter, and Mark J. van der Laan

## Abstract

The conditional relative risk is an important measure in medical and epidemiological studies when the outcome of interest is binary (i.e. disease vs. no disease). When the outcome is common, estimation of conditional relative risk and related parameters can be problematic, especially when the exposure or covariates are continuous. We propose a new estimation procedure based on targeted maximum likelihood methodology that targets the parameters relating to the conditional relative risk for common outcomes under a log-linear, or multiplicative, semi-parametric model. In this paper, we present three possible targeted maximum likelihood estimators for relative risk parameters implied by such a model: log-binomial based, Poisson-based, and a general semi-parametric approach. We present the properties and trade-offs of each of these estimators, focusing in particular on the Poisson-based estimator, which is most practical for implementation. We show that the resulting estimator is double robust and asymptotically linear, and inference can be obtained using the corresponding influence curve. The robustness of our estimator is compared to alternative methods (e.g. log-linear, Poisson regression) through simulation under model misspecification and increasing violations of the positivity assumption. The estimation procedure is then applied to a study of HIV genetic susceptibility scores, which aims to determine the effects of different genetic susceptibility scores on viral response. Effect modi?cation by other covariates in the model is also explored.

# 1  Introduction

In medical and epidemiological studies, when an outcome of interest is binary (i.e. disease vs. no disease) researchers are often interested in measuring the relative risk of an exposure or treatment. The conditional relative risk (RR), which adjusts for possible confounders, can be an important measure in such studies, because of its easy interpretation as the increase in the probability of the outcome in an exposed group relative to an non-exposed group. This straightforward interpretation allows the researcher to communicate results to a variety of audiences. Given a rare outcome, the conditional RR can be approximated by the conditional odds ratio (OR) and easily estimated using logistic regression. However, when the outcome is common, estimating the conditional RR can be problematic, especially when the exposure or covariates are continuous(McNutt et al., 2003; Barros and Hirakata, 2003; Lumley et al., 2006). In this paper, we propose a new approach for estimating the conditional RR and its related parameters for a common outcome based on targeted maximum likelihood methodology (van der Laan and Rubin, 2006). This approach is developed under a flexible multiplicative semi-parametric model and can be implemented using standard statistical software.

In the literature, the most prevalent methods for estimating conditional relative risk parameters for a common outcome are parametric methods based on log-linear (e.g. Wacholder (1986), Skov et al. (1998) and Zocchett et al. (1995)), Poisson (e.g. Zou (2004) and McNutt et al. (2003)), and Cox regression models (e.g. Lee and Chia (1993) and Lee (1994)). Overviews and comparisons of these methods can be found in papers by McNutt et al. (2003), Barros and Hirakata (2003) and Lumley et al. (2006). These three methods do not imply three unique estimators, however. As the above mentioned papers point out, the estimators implied by using Poisson regression and Cox regression are equivalent. In practice, Poisson regression is often preferred and easier to implement, therefore we focus our comparisons on Poisson regression. Some have also suggested methods to convert the OR to the RR. However, these methods are susceptible to bias in both their estimates, tests, and confidence intervals (Localio et al., 2007; McNutt et al., 2003).

If the regression model is correctly specified, these parametric methods will provide consistent estimates of conditional relative risk parameters. Log-binomial regression estimates the conditional RR directly using maximum likelihood estimation. However, when continuous covariates are included in the model, log-binomial regression becomes highly susceptible to convergence issues and requires modifications to achieve more stable parameter estimates (e.g. Wacholder (1986), SAS (2003) and Lumley et al. (2006)). Poisson regression is not plagued by these convergence issues and provides consistent estimates of the conditional relative risk when the model is correct (Stijnen and Van Houwelingen, 1993; Zou, 2004; Carter et al., 2005). However, the standard errors provided by Poisson regression are overestimated, and alternative methods such as the sandwich estimator (Zou, 2004; Carter et al., 2005) or bootstrap (Barros and Hirakata, 2003) must be used to obtain correct inference.

Due to their dependence on the accuracy of a fully specified model, parameter estimates from the above parametric methods are often biased in observational studies. More flexible alternatives presented in the literature include semi-parametric counterparts to the log-linear and Poisson regression models. Typically these are built under either generalized partial linear models or generalized additive partial linear models, the latter being a less flexible approach where each covariate has a separate additive component in the model (Hastie and Tibshirani, 1990). Estimation methods for the parameter under a partial linear model include profile likelihood methods (Speckman, 1988; Severini and Wong, 1992; Severini and Staniswalis, 1994) and backfitting algorithms (Hastie and Tibshirani, 1990). Under the additive partial linear model, estimation methods include backfitting (Buja et al., 1989) and marginal integration (Chen et al., 1996). Though these estimation methods are less dependent on model specification, they do not target estimation towards the parameter of interest and can still result in biased estimates and improper inference.

In this paper, we introduce double robust (DR) targeted maximum likelihood estimators (TMLEs) of conditional relative risk parameters. These estimators are developed under a multiplicative semi-parametric regression model, which is also referred to as a log-link generalized partial linear model. The model only requires specification of the model terms relating to the variable of interest (i.e. the exposure) and allows the remaining terms to be estimated as flexibly as possible, using a data adaptive approach. Targeted maximum likelihood estimation updates an initial estimator of the density of the data under the assumed model in

a direction that targets estimation towards the parameter of interest (van der Laan and Rubin, 2006). In this case, the parameter of interest is the coefficient (or coefficient vector) for the exposure variable (i.e. variable of interest) and any effect modifier (i.e. interaction) terms. The TMLE for these coefficients can be used to compute conditional relative risk. Previous applications of target maximum likelihood under a semi-parametric regression model can be found in Tuglus and van der Laan (2008, 2010), and additional applications of targeted maximum likelihood can be found in van der Laan et al. (September, 2009)

In this paper, we present three TMLE's for the parameter of interest in the multiplicative semi-parametric model, but focus in particular on one "practical" TMLE, which we recommend for use in practice. Each TMLE is based on the update of a different initial density and makes different assumptions on the distribution of the outcome. The three distributions for the TMLEs we present here are: (1) log-binomial, (2) Poisson, and (3) overdispersed exponential. The first TMLE, based on an update of a log-binomial density, is the natural TMLE for parameters related to the conditional relative risk. It correctly assumes a binary outcome and the resulting TMLE is double robust, asymptotically linear, and efficient. However, estimation requires log-binomial regression, which is often computationally unstable; therefore this TMLE is not typically feasible in practice. The second TMLE is based on an update of a Poisson density. This TMLE assumes the outcome follows a Poisson distribution, which is incorrect for a binary outcome; therefore, given a binary outcome, the resulting estimator is double robust and asymptotically linear but no longer efficient. However, this TMLE is computationally stable and straightforward to implement with correct inference. Therefore, this TMLE is considered the most practical of the three estimators, and this paper will focus primarily on this TMLE for implementation and application. The third TMLE makes no assumptions on the distribution of the outcome and is based on an update of a density in the overdispersed exponential family. This TMLE is the most general of the three, and the targeted maximum likelihood updates of the two prior methods can be derived directly from the update for this TMLE. However, this method is not straightforward to implement in practice using standard software.

The layout of the paper is as follows. In Section 2, we present the data structure. In Section 3, we present the assumed multiplicative semi-parametric model, and in Section 4 we formalize the parameter of interest, which is implied by the model. In Section 5, we present the three TMLE's mentioned above, followed by step-by-step implementation instructions for the more practical Poisson-based TMLE in Section 6. In Section 7, we demonstrate the performance of this Poisson-derived TMLE with results from a variety of simulations, and in Section 8, we apply this TMLE in an HIV viral response data set. We conclude by summarizing the value our new approach in the discussion in Section 9.

## 2  Data Structure

Consider an observed point treatment data set consisting of $n$ independent and identically distributed (i.i.d.) observations of $O = (W, A, Y) \sim P_0 \in \mathcal{M}$. $W$ is a vector of baseline covariates, $A$ is the exposure of interest, and $Y = \{0, 1\}$ is a binary outcome. $P_0$ denotes the true distribution of $O$, from which all subjects are sampled. $P_0$ is an element of a statistical model $\mathcal{M}$, which is a semi-parametric model defined below in Section 3.

In this article, the subscript 0 (e.g. $\beta_0$) will represent a parameter or function under the observed data generating distribution, and the subscript $n$ (e.g. $\beta_n$) will represent a estimate or estimator of a parameter or function determined from the sample population. Additionally a superscript indicates the iteration value, where a superscript 0 designates an initial value, a superscript $k$ designates the values at the $k^{th}$ iteration, and a superscript $*$ indicates the final converged value.

Note that for causal effects, we assume that $O$ is a missing data structure on a hypothetical full data structure $X = (W, Y_a : a \in A)$, which contains all counterfactual outcomes $Y_a$. We therefore view $A$ as the missingness variable, as $O$ contains only one of all possible counterfactual outcomes, $Y = Y_A$.

# 3    Multiplicative Semi-parametric Model

The likelihood of the observed data can be factorized as follows in terms of the observed data structure defined above:

$$P_0(O) = P_0(W)P_0(A|W)P_0(Y|A,W).$$

We make no assumptions about the distributions of $P_0(W)$ or $g_0(A|W) = P_0(A|W)$ and assume the following semi-parametric multiplicative model for the mean of $P_0(Y|A,W)$:

$$P_0(Y = 1|A,W) = e^{m_{\beta_0}(A,V)}P_0(Y|A = 0, W),$$

or

$$\log(P_0(Y = 1|A,W)) = m_{\beta_0}(A,V) + \log(P_0(Y = 1|A = 0, W)),$$

where $m_{\beta_0}(A,V)$ is a specified function of $A$ and effect modifiers $V \subset W$, and $P_0(Y = 1|A = 0, W)$ is unspecified. The model $m_{\beta_0}(A,V)$, can be any form such that $m_{\beta_0}(A = 0, V) = 0$ for all values, $v \in V$. We generally specify a linear function of $A$ and work with the following two models (1) simple main effect: $m_{\beta_0}(A,V) = \beta_0 A$ or (2) V-modified $m_{\beta_0}(A,V) = \beta_{0(1)}A + \beta_{0(2)}A : V$, where the effect of exposure $A$ is modified by covariate $V$. The latter model form is represented in terms of the parameter vector $\beta_0$ as $m_{\beta_0}(A,V) = \beta_0^T[A \ A : V]$.

Going forward, we define $\bar{Q}_0(A,W) \equiv P_0(Y = 1|A,W)$. For convenience, we introduce separate notation for $P_0(Y = 1|0, W)$ and define $\theta_0(W) \equiv \bar{Q}_0(0, W)$. We can then rewrite the multiplicative semi-parametric model as:

$$\bar{Q}_0(A,W) = e^{m_{\beta_0}(A,V)}\theta_0(W). \tag{1}$$

# 4    Parameter of Interest

We can represent conditional relative risk (RR) in terms of the observed data and the semi-parametric model as follows:

$$
\begin{aligned}
RR_0(a) \quad &= \frac{P_0(Y=1|A=a,W)}{P_0(Y=1|A=0,W)} \\
&= \frac{\bar{Q}_0(a,W)}{\theta_0(W)} \\
&= e^{m_{\beta_0}(a,V)},
\end{aligned}
$$

or on the log scale:

$$\log\left\{ \frac{\bar{Q}_0(a,W)}{\theta_0(W)} \right\} = m_{\beta_0}(a,V).$$

The model in 2 requires that both $\bar{Q}_0(A,W) > 0$ and $\theta_0(W) > 0$. In order for this to be true, we must have that $P_0(A = a|W) > 0$ for all $a, w$. This latter requirement is often referred to as the positivity assumption (Robins, 1986, 1987, 1999), or the experimental treatment assignment (ETA) assumption (Neugebauer and van der Laan, 2005). Therefore, the model in 2 is true only for $a, w$ for which there is support in the data. So to be more precise, we can write:

$$\frac{\bar{Q}_0(a,w)}{\theta_0(w)}I(a, w \in \mathcal{A}') = e^{m_{\beta_0}(a,v)}, \tag{2}$$

where $\mathcal{A}'$ contains the subset of $a, w$ for which the positivity assumption is not violated (i.e. for which $P_0(A = a|W) > 0$). This allows us to estimate the importance (i.e. risk) of exposure $A$ in predicting the

3

outcome $Y$, conditional on $W$, for those strata of $W$ in which the data have sufficient support. The analyst may want to allow for extrapolation across *all* $a, w$, but this would be unwise if the model was not true for across all $a, w$.

Based on the parameterization of the conditional relative risk shown above, we can define our parameter of interest as $\beta_0 = \Psi(P_0)$. We note that under the model defined in Section 3, the parameter of interest is only a function of $\bar{Q}_0(A, W) = P_0(Y = 1|A, W)$. Therefore we can denote, $\beta_0 = \Psi(\bar{Q}_0)$. If we define $m_{\beta_0}(A, V) = \beta_0 A$, then the parameter of interest, $\beta_0$ is the change in the log of the conditional relative risk for a one unit increment change in $A$. This allows us to estimate the importance (i.e. risk) of exposure $A$ in predicting the outcome $Y$, conditional on $W$, for those strata of $W$ in which the data have sufficient support.

In order for the parameter of interest to have a causal interpretation, we require that not only that $O$ is a missing data structure on a hypothetical full data structure $X = (W, Y_0, Y_1)$, as described in Section 2, but we also require the randomization assumption: $\{A \perp Y_0, Y_1|W\}$. However, even in the case where these assumptions do not hold, the variable importance parameter defined above is a well-defined and meaningful parameter.

We note that one motivation for using our multiplicative semi-parametric model is that regardless of whether or not we believe our model, we can still accurately test the following strong null hypothesis:

$$H_0 : \frac{P_0(Y = 1|A, W)}{P_0(Y = 1|A = 0, W)} = 1,$$

for all $W$. Under this null, the model is always correct. Therefore, we can construct valid tests of this null hypotheses.

# 5   Targeted Maximum Likelihood Estimation

Targeted maximum likelihood estimation updates an initial estimator of the density in a direction that targets estimation towards the parameter of interest. The update is achieved by regressing the outcome on a "clever covariate" while setting the initial estimate as an offset. This update is iterated until convergence. The "clever covariate" is derived such that the converged TMLE is also the solution to the double robust estimating equation for the parameter of interest. The converged TMLE is an asymptotically linear estimator of the parameter of interest, therefore inference can be based on the corresponding influence curve (van der Laan and Rubin, 2006).

For the parameter of interest, $\beta_0$, under the assumed semi-parametric multiplicative model defined in Section 3, three different TMLE's may be developed: (1) as an update to log-binomial density assuming the distribution for $P_0(Y|A, W)$ is binomial; (2) as an update to a Poisson density assuming the distribution for $P_0(Y|A, W)$ is Poisson; or (3) as an update to a general density of the overdispersed exponential family independent of distributional assumptions on $P_0(Y|A, W)$. The latter TMLE is the most general and assumes only a semi-parametric multiplicative model of general form as presented in 1. We note that TMLE's (1) and (2) are examples of TMLE (3) and can be derived from method (3) under their respective distributional assumption (see Appendix A for details).

In the case of a binary outcome, where one is interested in estimating the conditional relative risk, all three TMLE's are double robust and asymptotically linear estimators. However, only TMLE's (1) and (3) respect the binary form of the outcome and are therefore efficient. The Poisson-based TMLE (2) assumes the outcome follows a Poisson distribution. Therefore, when the outcome is binary, the initial estimator for TMLE (2) is always misspecified, and the corresponding influence curve is not the efficient influence curve. Although TMLE (2) is not efficient it remains DR - that is, the efficient score estimating function in the semi-parametric Poisson regression model is an unbiased DR estimating function for the parameter of interest in the semi-parametric conditional mean model, which does not assume a Poisson distribution. Consequently, the Poisson-based TMLE (2) is consistent, and we can provide correct inference using the corresponding influence curve.

In practice, the Poisson-based TMLE (2) is more stable than the TMLE (1), which requires updating a log-binomial regression. As mentioned previously, log-binomial regression is often unstable especially

4

given continuous covariates. Although TMLE (3) is more general with no distributional assumptions on $P_0(Y|A, W)$, it can not be implemented using standard software and is therefore a less practical estimator. Therefore in practice, the Poisson-based TMLE (2) is preferred and will be applied in this paper in simulation and application.

## 5.1 Specifics of the TMLE's

For each of the three TMLEs we present the form of the density, its fluctuation model, and the derived "clever covariate." In Appendix B we provide their respective efficient scores and estimating functions as well as the derivation of their respective "clever covariates."

For all three TMLEs, the initial density $P^0(Y|A, W)$, with mean, $E^0[Y|A, W] = P^0(Y = 1|A, W) = \bar{Q}^0(A, W)$ is defined in terms of the semi-parametric model presented above, where

$$\log \bar{Q}^0(A, W) = m_{\beta^0}(A, V) + \log \theta^0(W).$$

A class of submodels fluctuated with parameter $\epsilon$ is defined as $P^0(\epsilon)(Y|A, W)$ with corresponding mean

$$\log \bar{Q}^0(\epsilon)(A, W) = m_{\beta^0(\epsilon)}(A, V) + \log \theta^0(\epsilon)(W),$$

where $m_{\beta^0(\epsilon)}(A, V) = m_{\beta^0+\epsilon}(A, V)$ and $\log \theta^0(\epsilon)(W) = \log \theta^0(W) + \epsilon r_{\bar{Q}^0,g^0}(W)$. The form of $r_{\bar{Q},g}(W)$ is determined such that at $\epsilon = 0$, $\bar{Q}^0(\epsilon = 0)(A, W) = \bar{Q}^0(A, W)$ and the linear span of the score of the likelihood for $\bar{Q}^0(\epsilon)(A, W)$ with respect to $\epsilon$ at $\epsilon = 0$ includes the efficient score (van der Laan and Rubin, 2006).

Given a model $m_\beta(A, V)$ that is linear in $\beta$, the model,

$$\log \bar{Q}^0(\epsilon)(A, W) = m_{\beta^0(\epsilon)}(A, V) + \log \theta^0(\epsilon)(W),$$

can be rearranged as an update to the initial fit

$$\log \bar{Q}^0(\epsilon)(A, W) = \log \bar{Q}^0(A, W) + \epsilon^T \frac{d}{d\beta^0} m_{\beta^0}(A, V) + \epsilon^T r_{\bar{Q}^0,g^0}(W).$$

The update can be achieved by estimating $\epsilon$ with standard maximum likelihood estimation. For log-binomial and Poisson methods, the update can be completed using generalized linear regression setting the initial estimate, $\bar{Q}^0(A, W)$, as an offset and regressing $Y$ onto the following "clever covariate",

$$H^*_{\bar{Q}^0,g^0}(A, W) = \frac{d}{d\beta^0} m_{\beta^0}(A, V) + r_{\bar{Q}^0,g^0}(W).$$

In practice the update process is iterated such that $\beta^{k+1} = \beta^k + \epsilon^k$ and $\log \theta^{k+1}(W) = \log \theta^k(W) + \epsilon^k r_{\bar{Q}^k,g^k}(W)$, where $\epsilon^k$ is the update parameter for the $k^{th}$ update (e.g. $\beta^1 = \beta^0 + \epsilon^0$). Convergence is achieved when $\epsilon^k \approx 0$. The final converged estimate is the solution to the robust estimating equation corresponding to the efficient score

$$\frac{1}{n} \sum_{i=1}^{n} \left[ D_{h^*,\bar{Q}^*,g}(O_i) \right] = 0,$$

Therefore, the TMLE, $\beta^*$, inherits the DR properties of the solution to the efficient estimating equation and the efficient influence curve can be used to estimate the correct covariance and inference.

### 5.1.1 TMLE (1): Log-binomial model

Assuming a log-binomial distribution, the initial density is a binomial density defined as

$$P^0(Y|A, W) = \bar{Q}^0(A, W)^Y (1 - \bar{Q}^0(A, W))^{1-Y},$$

where $\bar{Q}^0(A, W) = P^0(Y = 1|A, W) = \theta^0(W) e^{m_{\beta^0}(A,V)}$ with the associated fluctuation

$$P^0(\epsilon)(Y|A, W) = \bar{Q}^0(\epsilon)(A, W)^Y (1 - \bar{Q}^0(\epsilon)(A, W))^{1-Y},$$

5

given $\bar{Q}^0(\epsilon)(A,W) = \theta^0(\epsilon)(W)e^{m_{\beta^0(\epsilon)}(A,V)}$.

The proper form of the fluctuation function $r_{\bar{Q}^0,g^0}(W)$ is defined as follows

$$r_{\bar{Q}^0,g^0}(W) = -\frac{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\frac{d}{d\beta^0}m_{\beta^0}(A,V)\Big|W\right]}{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\Big|W\right]},$$

The update is completed using log-binomial regression with an offset equal to the initial fit and a "clever covariate" defined as

$$H^*_{\bar{Q}^0,g^0}(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\frac{d}{d\beta^0}m_{\beta^0}(A,V)\Big|W\right]}{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\Big|W\right]}.$$

### 5.1.2 TMLE (2): Poisson

Under the Poisson distribution, the initial density is a Poisson density defined as

$$P^0(Y|A,W) = \frac{\bar{Q}^0(A,W)^Y}{Y!}e^{-\bar{Q}^0(A,W)},$$

with the associated fluctuation

$$P^0(\epsilon)(Y|A,W) = \frac{\bar{Q}^0(\epsilon)(A,W)^Y}{Y!}e^{-\bar{Q}^0(\epsilon)(A,W)}.$$

The proper form for $r_{\bar{Q}^0,g^0}(W)$ is

$$r_{\bar{Q}^0,g^0}(W) = -\left\{\frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E[\frac{d}{d\beta^0}m_{\beta^0}(A,V)e^{m_{\beta^0}(A,V)}|W]}{E[e^{m_{\beta^0}(A,V)}|W]}\right\}.$$

The update is completed using Poisson regression with an offset equal to the initial fit and "clever covariate" defined as

$$H^*_{\bar{Q}^0,g^0}(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \left\{\frac{E[\frac{d}{d\beta^0}m_{\beta^0}(A,V)e^{m_{\beta^0}(A,V)}|W]}{E[e^{m_{\beta^0}(A,V)}|W]}\right\}.$$

In practice, the Poisson based TMLE is recommended due to its computational stability. Therefore, subsequent implementation instructions and applications in this chapter will be presented in terms of the Poisson TMLE.

### 5.1.3 TMLE (3): Overdispersed Exponential Family

A density of the overdispersed exponential family in canonical form is represented as follows (McCullagh and Nelder, 1989)

$$P^0(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\eta^0 Y - B(\eta^0)}{d(\tau)}\right\}.$$

for this subclass of models $d(\tau)$ is the conditional residual variance, $\sigma_Y^2(A,W)$, $\eta^0 = \bar{Q}^0(A,W) = \theta^0(W)\exp(m_{\beta^0}(A,V))$ and $B'(\eta^0) = \bar{Q}^0(A,W)$. Therefore it can be rewritten as follows

$$P^0(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\bar{Q}^0(A,W)Y - B(\bar{Q}^0(A,W))}{\sigma_Y^2(A,W)}\right\}.$$

We define a class of submodels, fluctuated by parameter $\epsilon$ as

$$P^0(\epsilon)(Y|A,W) = h_e(Y,\tau)\exp\left\{\frac{\bar{Q}^0(\epsilon)(A,W)Y - B(\bar{Q}^0(\epsilon)(A,W))}{\sigma_Y^2(A,W)}\right\},$$

6

where $\bar{Q}^0(\epsilon)(A,W) = \theta(\epsilon)(W)\exp(m_{\beta(\epsilon)}(A,V))$, $\theta(\epsilon)(W) = \theta(W)\exp(\epsilon r_{\bar{Q}^0,g^0}(W))$ and $\beta(\epsilon) = \beta + \epsilon$.

The proper form of the fluctuation function, $r_{\bar{Q}^0,g^0}(W)$, is defined as follows

$$r_{\bar{Q}^0,g^0}(W) = -\frac{E\left[\frac{d}{d\beta^0}m_{\beta^0}(A,V)\frac{Q^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E\left[\frac{Q^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}$$

and the clever covariate is defined as

$$H_{\bar{Q}^0,g^0}^*(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E\left[\frac{d}{d\beta^0}m_{\beta_0}(A,V)\frac{Q^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E\left[\frac{Q^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}$$

## 5.2 Estimating the "clever covariate"

Unlike the TMLE for the additive semi-parametric regression model for a continuous outcome presented in Tuglus and van der Laan (2008, 2010), the "clever covariate" presented here is not only a function of the conditional mean of $A$, given $W$, but depends on $P(A\,|\,W)$ in a more complex way. When $A$ is continuous, the "clever covariate" can be directly estimated using a data-adaptive regression algorithm to estimate the expectations in the numerator and denominator of the second term. This is the method that is used in the following application in Section 8.

Given a binary or categorical $A$ with $L$ levels, the "clever covariate" can also be determined directly as follows

$$H_{\bar{Q},g}^*(A,W) = A - \frac{\sum_{l=1}^L \frac{d}{d\beta}m_\beta(A=a_l,V)e^{m_\beta(A=a_l,V)}P(A=a_l|W)}{\sum_{l=1}^L e^{m_\beta(A=a_l,V)}P(A=a_l|W)}.$$

This method requires the estimation of $P(A=a_l|W)$ for all $L$. It is recommended that these values are estimated data-adaptively. This method can also be used to approximate the expectation given a continuous $A$ if the analyst is willing to discretize $A$ solely for the purpose of estimating this covariate.

## 5.3 Inference

An important feature of targeted maximum likelihood estimators is that they solve the corresponding robust estimating equation as presented in Section 5. Therefore, statistical inference of the converged estimate of $\beta_0$ can be based on the influence curve associated with this estimating equation. Covariance estimates based on the influence curve are consistent if the estimator for $g_0(W) = P_0(A|W)$ is correctly specified and the estimator for $\bar{Q}_0$ is misspecified. If $g_n$ is correctly specified, the covariance is known to be asymptotically conservative (van der Laan et al., September, 2009), except if $\theta_n$ is consistent or if $g_n = g_0$ is known. Covariance may also be estimated using the bootstrap. However, this TMLE requires iteration and can be computationally intensive.

For a parameter vector $\beta_0$ of length $p$, an estimate of the $p$ by $p$ covariance matrix, $\Sigma_n$, can be obtained using the influence curve defined for a single subject as

$$IC(O) = c^{-1}D_{h_0,\bar{Q}_0,g_0}(O),$$

given scale factor

$$c = -\mathbb{E}\left[\frac{d}{d\beta_0}D_{h_0,\bar{Q}_0,g_0}(O)\right]$$

where $IC(O)$ is a 1 by $p$ vector for a parameter vector $\beta_0$ of length $p$.

The empirical estimate for the covariance of $\beta_n$ is

$$\Sigma_n = \frac{1}{n}\sum_i \widehat{IC}_n(O_i)\widehat{IC}_n(O_i)^T$$

7

and the normal approximation

$$\sqrt{n}(\beta_n^* - \beta_0) \sim N(0, \Sigma_n),$$

can be used for the purpose of statistical inference.

Using the estimated covariance matrix, hypothesis tests can be performed for a single component $\beta_n(j)$, where $j = 1, \ldots, p$, under the null hypothesis $H_0 : \beta_0(j) = 0$ using a standard test statistic to obtain p-values, with estimated variance $\Sigma_n(j, j)$:

$$T_n(j) = \frac{\sqrt{n}\beta_n(j)}{\sqrt{\Sigma_n(j,j)}} \underset{n \to \infty}{\sim} N(0, 1).$$

Likewise the hypothesis $H_0 : c^T \beta_0 = 0$ can also be tested using a standard Wald test, where the covariance of $c^T \beta_n$ is $c^T \Sigma_n c$. In the application presented in this paper, effect modification is tested at various levels of $V$. In this situation the vector $c$ for the Wald test is $\{1, V_q\}$, where $V_q$ is a specific quantile of effect modifier $V$.

# 6 Step-by-Step Implementation

In this section, we provide step-by-step instructions for implementing targeted maximum likelihood estimation for $\beta_0$ under the semi-parametric regression model presented in (1). Although the following steps are presented assuming a Poisson distribution, the TMLE under the log-binomial distribution as well as the more general TMLE based on the overdispersed exponential density can be implemented in a similar fashion by substituting in the appropriate clever covariate and using log-binomial regression for the update.

For the following implementation steps, we assume the general linear model form $m_\beta(A, V) = A(\beta^T V)$.

**(1) Obtain an initial estimate,** $\bar{Q}_n^0(A, W)$ respecting the semi-parametric form $\log\bar{Q}_n^0 = m_{\beta_n^0} + \log\theta_n^0$, which provides the initial estimate for the parameter of interest, $\beta_n^0$. This can be accomplished by fitting the semi-parametric model using methods such as those described in Speckman (1988); Severini and Wong (1992); Hastie and Tibshirani (1990), or by using methods such as DSA Sinisi and van der Laan (2004) to fix the parametric portion of the model and allow the rest to be estimated data-adaptively. A more flexible alternative, uses an initial estimate for $\bar{Q}_0(A, W)$ of general model form obtained using data-adaptive machine learning algorithms such as super learner van der Laan et al. (2007). We estimate $\bar{Q}_n(A = 0, W)$ using this general model fit and then regress $Y$ onto the model $m_{\beta^0}(A, W)$ and treat $\bar{Q}_n(A = 0, W)$ as a covariate. This results in our initial density estimate, $\bar{Q}^0(A, W)$ and initial estimate, $\beta^0$ under the correct model.

**(2) Calculate the "clever covariate"** , $H_{\bar{Q}_n^0, g_n}^*(A, W) = AV - \frac{E[AVe^{(A(\beta_n^0{}^T V)}|W]}{E[e^{(A(\beta_n^0{}^T V))}|W]}$. If $A$ is discrete the $H_{\bar{Q}_n^0, g_n}^*$ can be calculated directly as described in Section 5.2, or by estimating the numerator and denominator of the second term ($E[AVe^{(A(\beta_n^0{}^T V))}) | W]$ and $E[e^{(A(\beta_n^0{}^T V))}) | W]$ respectively) using a flexible data-adaptive algorithm.

**(3) Estimate the fluctuation parameter,** $\epsilon^0$ using Poisson regression to project $Y$ onto $H_{\bar{Q}_n^0, g_n}^*$, while setting the initial fitted values, $\bar{Q}_n^0(A, W)$, as an offset. The estimated coefficient associated with the "clever covariate" is the fluctuation parameter estimate, $\epsilon_n^0$. If completed in R this is equivalent to fitting the following regression formula $Y \sim H_{\bar{Q}_n^0, g_n}^*(A, W) + offset(\bar{Q}_n^0(A, W)) - 1$ using `glm`. Note that there is no intercept, only the offset value.

**(4) Update initial estimate** by setting $\beta_n^1 = \beta_n^0 + \epsilon_n^0$ and and setting the updated fit as $\log(\bar{Q}_n^1(A, W)) = \log(\bar{Q}_n^0(A, W)) + \epsilon_n^0 H_{\bar{Q}_n^0, g_n}^*(A, W)$. These are the first-step updates.

**(5) Iterate** steps 1 through 4. At each $k^{th}$ iteration, set $\beta_n^k = \beta_n^{k-1} + \epsilon_n^{k-1}$ and $\log(\bar{Q}_n^k(A, W)) = \log(\bar{Q}_n^{k-1}(A, W)) + \epsilon_n^{k-1} H_{\bar{Q}_n^0, g_n}^*(A, W)$, until convergence ($\epsilon_n^k \approx 0$).

8

**(6) Obtain standard error and inference** for the final converged estimate $\beta_n^*$ using the influence curve or bootstrap estimates, and then calculate standard errors, p-values and confidence intervals as described in Section 5.3.

# 7 Simulations

In this section, we assess the properties of the Poisson-derived TMLE of $\beta_0$ with simulations that cover a range of scenarios seen in actual data sets. With these simulations, we demonstrate the double robustness properties of this TMLE and differences in variability when we estimate $\bar{Q}_0(A, W)$ and/or $g_0(A|W)$ consistently or with super learner. We compare our results to those from common estimators in the literature (those obtained from parametric log-binomial and Poisson regression models), as well as to the estimator obtained from an original fit of the semi-parametric model, using super learner to estimate $\theta_0$. To evaluate the performance of all estimators, we focus on bias, variance, mean squared error and confidence intervals.

We show that the relative performance of the TMLE, when compared to the other estimators, depends partly on the level of sparsity in the data. As discussed earlier, within strata defined by $W$, we would like the probability of exposure $A$ to be bounded away from zero and one. When $A$ is binary or categorical and is perfectly randomized, this is guaranteed. The common methods in the literature perform well under this scenario. However, when $A$ is continuous and in observational studies, analysts often have the challenge that $g_0(A|W)$ is very small for some strata of $W$. In other words, we have practical violations of the positivity assumption for all $a, W$.

Different estimators are affected by positivity violations in different ways. In the following simulations, we demonstrate the relative performance the Poison-derived TMLE of $\beta_0$, for both binary and continuous $A$, (1) when $A$ is perfectly randomized, (2) when the relationship between $A$ and $Y$ is confounded by $W$ but there are no positivity violations and (3) when there are extreme positivity violations.

In all simulations, the outcome variable, $Y \in \{0, 1\}$ is an indicator, such as of of disease status. Also in all of the simulations, $W$ is a vector of five covariates, which were generated as follows:

$$
\begin{aligned}
W_1 &\sim Binom(1, 0.3) \\
W_2 &\sim Binom(1, 0.65) \\
W_3 &\sim N(0, 2) \\
W_4 &\sim N(100, 10) \\
W_5 &\sim N(1, 0.3).
\end{aligned}
$$

We set
$$
\bar{Q}_0(A, W) = e^{-0.1A} e^{I + 0.1W_3 + 0.02W_2W_3 - 0.01W_1W_4 - 0.02W_5}, \tag{3}
$$
where $I$ takes the following values for the various simulations:

| | Binary A | | | Continuous A | |
|---|---|---|---|---|---|
| Simulation 1 | Simulation 2 | Simulation 3 | Simulation 1 | Simulation 2 | Simulation 3 |
| -0.8 | -0.8 | -1.0 | -0.8 | -0.4 | -1.4 |

As (3) shows, $\beta_0 = -0.1$ in all simulations. Also, (3) shows that $m_{\beta_0}(A) = \beta_0 A$, so we have assumed there are no effect modifiers.

## 7.1 Simulations for Binary A

For a binary $A$, we consider the following three conditional distributions for $g_0(A|W)$:

9

1. For the first simulation, $A$ is perfectly randomized such that $g_0(A|W) = 0.5$. When $g_0(A|W)$ is misspecified for this simulation, $g_n(A|W) = 0.6$.

2. For the second simulation, $A$ is dependent on $W$ such that:

$$g_0(A|W) = \frac{1}{1 + exp(-(0.1W_3))}.$$

With this mechanism for exposure, the correlation between $A$ and $W_3$ is 0.10, and values of $g_0(A|W)$ range from 0.28 to 0.73, with a median of 0.49. Therefore, we do not have positivity violations. For this simulation, when $g_n(A|W)$ is misspecified, the estimator depends only on $W_1$.

3. For the third simulation, $A$ is again dependent on $W$, but now we have:

$$g_0(A|W) = \frac{1}{1 + exp(-(1.0W_3))}.$$

This mechanism for exposure leads to positivity violations because $g_0(A|W) \in [5.4 \times 10^{-5}, 1.0]$. The median value is 0.54. The correlation between $A$ and $W_3$ is now 0.61. Misspecification of $g_n(A|W)$ again occurs by having the estimator only depend on $W_1$.

## 7.2   Simulations for Continuous A

For continuous $A$, we again varied $g_0(A|W)$ three ways:

1. For the first simulation, $A$ is not dependent on $W$. It is normally distributed such that $A \sim N(1, 0.6)$.

2. For the second simulation, $A$ is dependent on $W$ such that $A \sim N(1, 0.6) + 0.1W_3$. In this simulation, the correlation between $A$ and $W_3$ is $-0.1$.

3. For the third simulation, $A$ is dependent on $W$ such that $A \sim N(0, 0.6) - 0.8W_3$. The correlation between $A$ and $W_3$ in this simulation is $-0.6$.

For all simulations, we generated 1000 samples of size 1000. All data were generated and all estimators were implemented using R (Team, 2010).

## 7.3   Simulation Results

Tables 1.1 and 1.2 present results for estimating $\beta_0$ when $A$ is binary and when $A$ is continuous. For continuous $A$, we estimated the numerator and denominator of the clever covariate using the `lars` package in R (Efron et al., 2003). The first column of the tables presents the initial substitution estimator, $\beta_n^0$, based on the initial estimate of $\bar{Q}_0$. The second column presents the TMLE, $\beta_n^*$, obtained by substitution after $k$ iterations of updating $\bar{Q}_n^0$ to obtain $\bar{Q}_n^*$. The subsequent columns provide the bias, mean squared error (MSE) and empirical variance based on the estimates from 1000 samples. We also include the mean of the variance estimates calculated from the empirical variance of efficient influence curve divided by the sample size of 1000. Finally, the last column shows the coverage probability (CP), or the percentage of the time that the 95% confidence interval contains the true value of $\beta_0 = -0.1$.

Each panel in the tables corresponds to the simulations described above, and the rows in each panel indicate the specification of the estimators of $\bar{Q}_0$ and $g_0$, on which the TMLE is based. For example, "Qcgc" indicates that the correct terms were included when estimating both $\bar{Q}_0$ and $g_0$. "Qcgm" indicates that the

10

Table 1: Performance of Poisson-derived TMLE, binary A, by simulation

| | $\beta_n^0$ | $\beta_n^*$ | Bias | MSE | $\text{Var}(\beta_n^*)$ | $\text{Var}(IC_{eff})/n$ | CP |
|---|---|---|---|---|---|---|---|
| **Simulation 1** | | | | | | | |
| Qcgc | -0.102 | -0.102 | -0.002 | 0.006 | 0.006 | 0.007 | 0.964 |
| Qcgw | -0.102 | -0.102 | -0.002 | 0.006 | 0.006 | 0.011 | 0.990 |
| Qwgc | -0.101 | -0.101 | -0.001 | 0.006 | 0.006 | 0.007 | 0.966 |
| Qslgsl | -0.090 | -0.102 | -0.002 | 0.006 | 0.006 | 0.007 | 0.960 |
| **Simulation 2** | | | | | | | |
| Qcgc | -0.103 | -0.103 | -0.003 | 0.007 | 0.007 | 0.007 | 0.950 |
| Qcgw | -0.103 | -0.103 | -0.003 | 0.007 | 0.007 | 0.007 | 0.952 |
| Qwgc | -0.057 | -0.102 | -0.002 | 0.007 | 0.007 | 0.007 | 0.944 |
| Qslgsl | -0.088 | -0.101 | -0.001 | 0.007 | 0.007 | 0.007 | 0.948 |
| **Simulation 3** | | | | | | | |
| Qcgc | -0.111 | -0.111 | -0.011 | 0.017 | 0.017 | 0.016 | 0.950 |
| Qcgw | -0.111 | -0.111 | -0.011 | 0.016 | 0.016 | 0.008 | 0.816 |
| Qwgc | 0.169 | -0.109 | -0.009 | 0.017 | 0.017 | 0.016 | 0.944 |
| Qslgsl | -0.091 | -0.109 | -0.009 | 0.017 | 0.017 | 0.016 | 0.940 |

estimator of $g_0$ was misspecified as described above, while the estimator for $\bar{Q}_0$ included the correct terms; and "Qgmc" indicates that the estimator for $\bar{Q}_0$ was misspecified as described above, while the correct terms were included when estimating $g_0$. Finally "Qslgsl" indicates that the super learner was used for the estimators for both $\bar{Q}_0$ and $g_0$.

Tables 1.1 and 1.2 illustrate the properties we expect to see for the TMLE of $\beta_0$:

- The TMLE is double-robust. The finite-sample bias is close to zero if the estimator for either $\bar{Q}_0$ or $g_0$ is consistent. We achieve this even under substantial confounding and extreme violations of positivity in Simulation 3.

- When the estimator for $g_0$ is consistent, the variance estimate obtained from the empirical variance of the efficient influence curve is approximately equal to the variance of the 1000 TMLEs and the coverage probability is approximately 95%. When the estimator for $g_n$ is inconsistent, this variance estimate is asymptotically conservative.

- Using the super learner to estimate both $\bar{Q}_0$ and $g_0$ provides robust estimators of either $\bar{Q}_0$ or $g_0$ so that we achieve comparable bias and variance as obtained when correctly specifying the models for $\bar{Q}_0$ and/or $g_0$. The one exception is when $A$ is continuous and we have extreme positivity violations.

Tables 1.3 and 1.4 compare the performance of the TMLE of $\beta_0$ to the common estimators in the literature when the initial working model for $\bar{Q}_0(A, W)$ is *incorrect*. All of the estimators in the literature will perform well when the parametric models on which they rely are correctly specified. However, we are very doubtful that anyone can ever specify a parametric model correctly. Therefore, we present comparisons under a more realistic scenario.

Within each panel for each simulation, the first two rows present results (bias, variance and MSE) for the common methods in the literature - using log binomial and Poisson regression to estimate W-adjusted relative risk. The third and fourth rows present results for the initial estimate of $\beta_0$, $\beta_n^0$, when $\bar{Q}_n^0(A, W)$ is incorrectly specified and when it is estimated by super learning. The last two rows then present results for the TMLE, $\beta_n^*$ when $\bar{Q}_n^0(A, W)$ is incorrectly specified and when it is estimated by super learning.

Figure 1.1 also compares the performance of TMLE to other relative risk estimators The following summarizes key observations from both the tables and figure:

Table 2: Performance of Poisson-derived TMLE, continuous A, by simulation

|  | $\beta_n^0$ | $\beta_n^*$ | Bias | MSE | $\text{Var}(\beta_n^*)$ | $\text{Var}(IC_{eff})/n$ | CP |
|---|---|---|---|---|---|---|---|
| **Simulation 1** | | | | | | | |
| Qcgc | -0.099 | -0.099 | 0.001 | 0.005 | 0.005 | 0.005 | 0.946 |
| Qcgw | -0.099 | -0.099 | 0.001 | 0.005 | 0.005 | 0.005 | 0.946 |
| Qwgc | -0.098 | -0.098 | 0.002 | 0.005 | 0.005 | 0.005 | 0.950 |
| Qslgsl | -0.089 | -0.099 | 0.001 | 0.005 | 0.005 | 0.005 | 0.950 |
| **Simulation 2** | | | | | | | |
| Qcgc | -0.098 | -0.098 | 0.002 | 0.005 | 0.005 | 0.005 | 0.956 |
| Qcgw | -0.098 | -0.098 | 0.002 | 0.005 | 0.005 | 0.005 | 0.956 |
| Qwgc | 0.016 | -0.099 | 0.001 | 0.005 | 0.005 | 0.005 | 0.956 |
| Qslgsl | -0.089 | -0.099 | 0.001 | 0.005 | 0.005 | 0.005 | 0.958 |
| **Simulation 3** | | | | | | | |
| Qcgc | -0.102 | -0.103 | -0.003 | 0.003 | 0.003 | 0.065 | 0.956 |
| Qcgw | -0.102 | -0.103 | -0.003 | 0.003 | 0.003 | 0.065 | 0.956 |
| Qwgc | 0.025 | -0.103 | -0.003 | 0.003 | 0.003 | 0.004 | 0.958 |
| Qslgsl | -0.096 | -0.105 | -0.005 | 0.003 | 0.003 | 0.003 | 0.948 |

Table 3: Relative performance of Poisson-derived TMLE, binary A

|  | Bias | Var | MSE |
|---|---|---|---|
| **Simulation 1** | | | |
| Log Binomial, incorrect | -0.004 | 0.007 | 0.007 |
| Poisson, incorrect | -0.004 | 0.007 | 0.007 |
| $\beta_n^0$, incorrect | -0.001 | 0.006 | 0.006 |
| $\beta_n^0$, SL | 0.010 | 0.008 | 0.008 |
| $\beta_n^*$ Qwgc | -0.001 | 0.006 | 0.006 |
| $\beta_n^*$ Qslgsl | -0.002 | 0.006 | 0.006 |
| **Simulation 2** | | | |
| Log Binomial, incorrect | 0.046 | 0.007 | 0.009 |
| Poisson, incorrect | 0.047 | 0.007 | 0.009 |
| $\beta_n^0$, incorrect | 0.043 | 0.007 | 0.009 |
| $\beta_n^0$, SL | 0.012 | 0.009 | 0.009 |
| $\beta_n^*$ Qwgc | -0.002 | 0.007 | 0.007 |
| $\beta_n^*$ Qslgsl | -0.001 | 0.007 | 0.007 |
| **Simulation 3** | | | |
| Log Binomial, incorrect | 0.277 | 0.010 | 0.087 |
| Poisson, incorrect | 0.277 | 0.010 | 0.087 |
| $\beta_n^0$, incorrect | 0.269 | 0.010 | 0.083 |
| $\beta_n^0$, SL | 0.009 | 0.016 | 0.016 |
| $\beta_n^*$ Qwgc | -0.009 | 0.017 | 0.017 |
| $\beta_n^*$ Qslgsl | -0.009 | 0.017 | 0.017 |

Table 4: Relative performance of Poisson-derived TMLE, continuous A

|  | Bias | Var | MSE |
|---|---|---|---|
| **Simulation 1** | | | |
| Log Binomial, incorrect | 0.000 | 0.004 | 0.004 |
| Poisson, incorrect | -0.002 | 0.005 | 0.005 |
| $\beta_n^0$, incorrect | 0.002 | 0.005 | 0.005 |
| $\beta_n^0$, SL | 0.011 | 0.007 | 0.007 |
| $\beta_n^*$ Qwgc | 0.002 | 0.005 | 0.005 |
| $\beta_n^*$ Qslgsl | 0.001 | 0.005 | 0.005 |
| **Simulation 2** | | | |
| Log Binomial, incorrect | 0.111 | 0.004 | 0.017 |
| Poisson, incorrect | 0.110 | 0.004 | 0.016 |
| $\beta_n^0$, incorrect | 0.116 | 0.005 | 0.018 |
| $\beta_n^0$, SL | 0.011 | 0.007 | 0.007 |
| $\beta_n^*$ Qwgc | 0.001 | 0.005 | 0.005 |
| $\beta_n^*$ Qslgsl | 0.001 | 0.005 | 0.005 |
| **Simulation 3** | | | |
| Log Binomial, incorrect | 0.120 | 0.000 | 0.015 |
| Poisson, incorrect | 0.123 | 0.000 | 0.015 |
| $\beta_n^0$, incorrect | 0.125 | 0.000 | 0.016 |
| $\beta_n^0$, SL | 0.004 | 0.003 | 0.003 |
| $\beta_n^*$ Qwgc | -0.003 | 0.003 | 0.003 |
| $\beta_n^*$ Qslgsl | -0.005 | 0.003 | 0.003 |

- In a randomized trial, as demonstrated in Simulation 1, all estimators perform comparably well, as expected, for both binary and continuous A.

- As the relationship between the true confounder and $A$ increases in Simulations 2 and 3, the estimators utilizing on log-binomial regression and Poisson regression are increasingly biased, while the variance (of the 1000 sample estimates of $\beta_0$) remains at the same or similar level (for binary A) or decreases (for continuous A).

- The TMLE of $\beta_0$ achieve the lowest MSE in both simulations with confounding (Simulations 2 and 3). We see a small trade-off in variance for removal of bias.

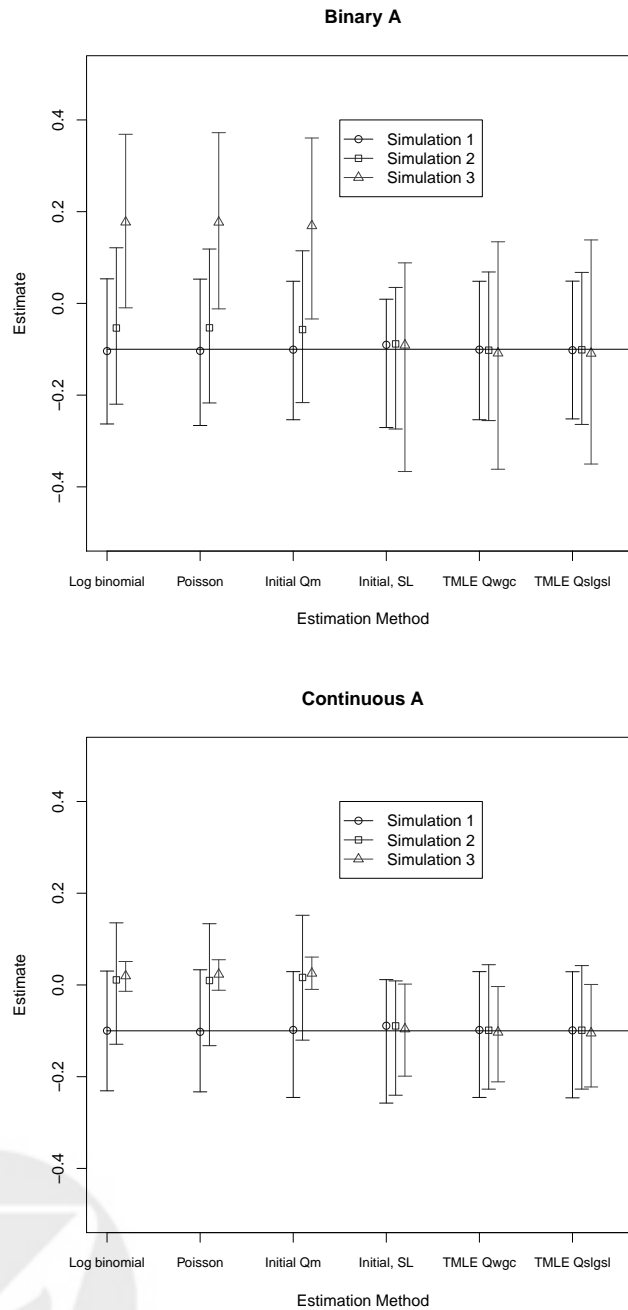- Even with positivity violations, the TMLE's are robust.

13

**Binary A**



**Continuous A**



Figure 1: Estimates and 95% confidence intervals by method

# 8 Application

Genotypic resistance testing has become a powerful tool for clinicians in determining the appropriate treatment regimen for people with HIV (Durant et al., 1999; Tural et al., 2002). However interpretation of the resistance testing can be difficult. Over the years multiple interpretation algorithms have been developed

to provide more straightforward measures of resistance Rhee et al. (2009). A study completed by Rhee et al. (2009) assessed the predictive ability of four genotypic resistance test interpretation algorithms (ANRS (de Recerche sur le SIDA , ANRS), HIVdb (Liu and Shafer, 2006), Rega (Van Laethem et al., 2002), and ViroSeq (Eshleman et al., 2004)) to determine virologic response (VR), adjusting for additional baseline covariates. In this analysis the data are reanalyzed using targeted maximum likelihood methodology to determine the importance of each algorithm with respect to its predictive fit. In other words, we are determining the importance of each genotypic algorithm ($A$) on VR ($Y$), adjusting for all other covariates in the original analysis ($W$). Then, the genotypic algorithm with the highest importance measure is analyzed further by estimating the modification of its effect by each of the covariates included in the original model.

## 8.1 Background

For individuals infected with human immunodeficiency virus (HIV) effective antiretroviral (ARV) treatments carry the promise of a longer and more gratifying life. HIV infects the body's immune system and progressively destroys and impairs its ability to fight off infection. ARV treatments are designed to slow down HIV reproduction and stall its debilitating effects. A properly designed and administered treatment can prolong survival and increase overall quality of life (of Health and on Clinical Practices for Treatment of HIV Infection A).

There are many ARV drugs available. Some of the more common classes of ARV drugs are boosted protease inhibitors (PIs), nucleoside reverse transcriptase inhibitors (NRTIs) and non-nucleoside reverse transcriptase inhibitors (NRTIs) (of Health and on Clinical Practices for Treatment of HIV Infection A). Patients are generally placed on a combination of several drugs from multiple classes called a treatment regimen. However, HIV is a rapidly evolving virus and commonly develops drug resistant mutations rendering initially effective treatment regimens useless (of Health and on Clinical Practices for Treatment of HIV Infection A). The rapid rate of mutation has made genotypic resistance testing essential to determining the appropriate treatment regimen for an individual patient (Durant et al., 1999; Tural et al., 2002). To facilitate the interpretation of these tests, interpretation algorithms have been developed. The algorithms are drug-specific and applied to a patients baseline genotype (Rhee et al., 2009).

## 8.2 Data

Study subjects were selected according to the eligibility requirements outlined in Rhee et al. (2009) from 16 clinics of the Kaiser-Permanente Medical Care Program, Northern California. In the original study 734 valid treatment change episodes (TCEs) were recorded for 641 patients. In this study a valid TCE occurs when a individual has undergone a change in treatment regimen within 24 weeks of a genotypic resistance test and has received at least four weeks of a new salvage regimen. Though the original study uses all TCEs, to simplify our analysis, we randomly select only one TCEs per patient.

Virologic response is measured according to plasma HIV-1 RNA levels. High plasma HIV-1 RNA levels indicate strong viral activity. The original study classifies VR into three classes: sustained, transient, and absent. Sustained VR is achieved when two subsequent tests show plasma RNA levels below the limit of quantification (BLQ). Transient refers to cases where only 1 subsequent test was at BLQ level, and absent when no subsequent test is BLQ. For this analysis sustained and transient classes are merged into a single class (VR=1) (see Rhee et al. (2009) for more details).

The original Rhee et al. (2009) study focused on four interpretation algorithms: ANRS (de Recerche sur le SIDA , ANRS), HIVdb (Liu and Shafer, 2006), Rega (Van Laethem et al., 2002), and ViroSeq (Eshleman et al., 2004). Each algorithm is applied to an individual patients baseline genotype to determine the drug specific genotypic susceptibility scores (GSSs) for each ARV. GSS measures range from 0 to 1, where a GSS of 1 indicates full susceptibility of HIV-1 to the particular ARV and a GSS of 0 indicates full resistance. The drug-specific GSSs are then combined into regimen specific GSSs (rGSSs) through three alternative weighing schemes: "boosted PI weighted", "comprehensive weighted", and "unweighted". The "unweighted" rGSS is the addition of all drug-specific GSSs weighted equally with 1.0. The "boosted PI weighted" rGSS increases the weight of the drug-specific GSSs for boosted PIs to 1.5. The "comprehensive weighted" rGSS increases

15

the weight of the drug-specific GSSs for boosted PIs to 2.0 and decreases the drug-specific GSSs for NRTIs to 0.5. This results in a total of 12 rGSSs. We standardize each rGSS by subtracting its mean and dividing by its standard deviation to allow direct comparison of the importance measures.

Additional covariates are also included in the original prediction analysis including individual demographics (age, sex, and race), features of ARV treatment prior to the TCE (i.e. duration of therapy, number of ARVs, etc.), and features of the salvage regimen (i.e. number of new ARVs, new ARV drug classes, etc.), as well as plasma HIV-1 RNA level and CD4 count at baseline.

## 8.3   Analysis

Targeted maximum likelihood estimation is first applied to estimate the variable importance of each rGSS with respect to its own prediction of VR. These estimates provide a measure of how much each variable changes the probability of VR on a relative scale. As stated previously in Section 5 , TMLE for variable importance updates an initial regression estimate to target the parameter of interest. In this case, the initial regression estimate for a specific rGSS is the super learner estimate, predicting VR using the rGSS and additional covariates. The covariate set is consistent across all rGSS fits and is defined using univariate logistic regression. Covariates associated with VR with a p-value of 0.1 or less are included. This analysis focuses on the overall effect of the rGSS, therefore the model is the simple effect model: $m_\beta(A, V) = \beta A$. For each importance measure, inference is obtained using the influence-curved based estimate of the standard error.

## 8.4   Results

Results of the initial analysis show significant importance values for all rGSS algorithms as expected. Note that though standardizing the scores allows us to directly compare the importance measures, the increment increase in RR is now relative to an increase of one in the z-score of the rGSS or correspondingly an increase of one standard deviation of the original rGSS measure. From Figures 2 and 3, we see that in general weighting did have an effect on overall importance. Weighting scheme 1 seems to increase the coefficient for each algorithm over the unweighted, and weighting scheme 2 increases it even more over that. Under any weighting scheme, ANRS seemed to have the highest coefficient and corresponding increment RR (e.g. the change in the probability of sustained virologic response for one unit increase in the zscore of the rGSS). Of the twelve, the highest importance is attributed to ANRS with "comprehensive weighting" with a coefficient $\beta = 0.206$ corresponding to an increment RR of 1.23 and associated p-value of 3.80e-06 before adjusting for multiple testing.

## 8.5   Secondary analysis

As a secondary analysis the rGSS with the highest importance, comprehensively weighted ANRS, is chosen and a V-modified variable importance analysis is performed, in which a covariate $V$ is included as an effect modifier of $A$. The model for this analysis is as follows

$$m_\beta(A, V) = A(\beta_A + \beta_v V)$$

where the parameter of interest is now measured as a function of two parameters with respect to a particular covariate, $V$. In this analysis, all other covariates are considered individually as effect modifiers for the selected rGSS, and the V-modified importance measure is estimated. The results are shown below.

The individual coefficient values estimated using TMLE with their respective confidence intervals are shown in Figure 4. Given only the coefficients it is difficult to interpret the results. To clarify, the increment RR change is calculated at varying levels, $V_q$ of each effect modifier $V$. This is defined as follows for value $V_q$ of any $V$ as

$$RR = e^{\beta_A + \beta_V V_q}.$$

16

Calculating the corresponding standard error is achieved by first calculating the standard error of the linear combination $c^T\beta = \beta_A + \beta_V V_q$, as $c\Sigma_V c^T$, where $\beta = \{\beta_A, \beta_V\}$, $c = \{1, V_q\}$, and $\Sigma_V$ is the influence curved based covariance estimate for $\{\beta_A, \beta_V\}$. Then, the delta method is used to calculate the corresponding standard error estimate for the exponential of this linear combination. The quantiles of $V$ (min, 25%, 50%, 75%, max) are chosen for $V_q$. Note that in some cases the covariate in binary and there are only two possible values, and therefore only two points. The results are shown in Figures 5 and 6.

From the results (Figures 5 and 6), it can be seen that the increase in the risk of sustained virologic response with respect to change in the rGSS score of ANRS under comprehensive weighting is modified by many of the covariates. This is not surprising due to the complexity of body's response to HIV. Virologic response is without a doubt a combination of genetics, current viral load, as well as current and past treatment regimens. For instance, it is logical that increased baseline viral load would result in increased risk of virologic response, but through this type of analysis, it can also be seen that increased baseline viral load seems to modify the effect of the genetic score on the relative risk of VR (Figures 5 and 6). This type of analysis helps elucidate and interpret the complex set of interactions that results in the body's virologic response. Having a method which targets the effect and provides consistent and locally efficient estimates of the effect with formal inference is key to further the understanding and treatment of diseases such as HIV.
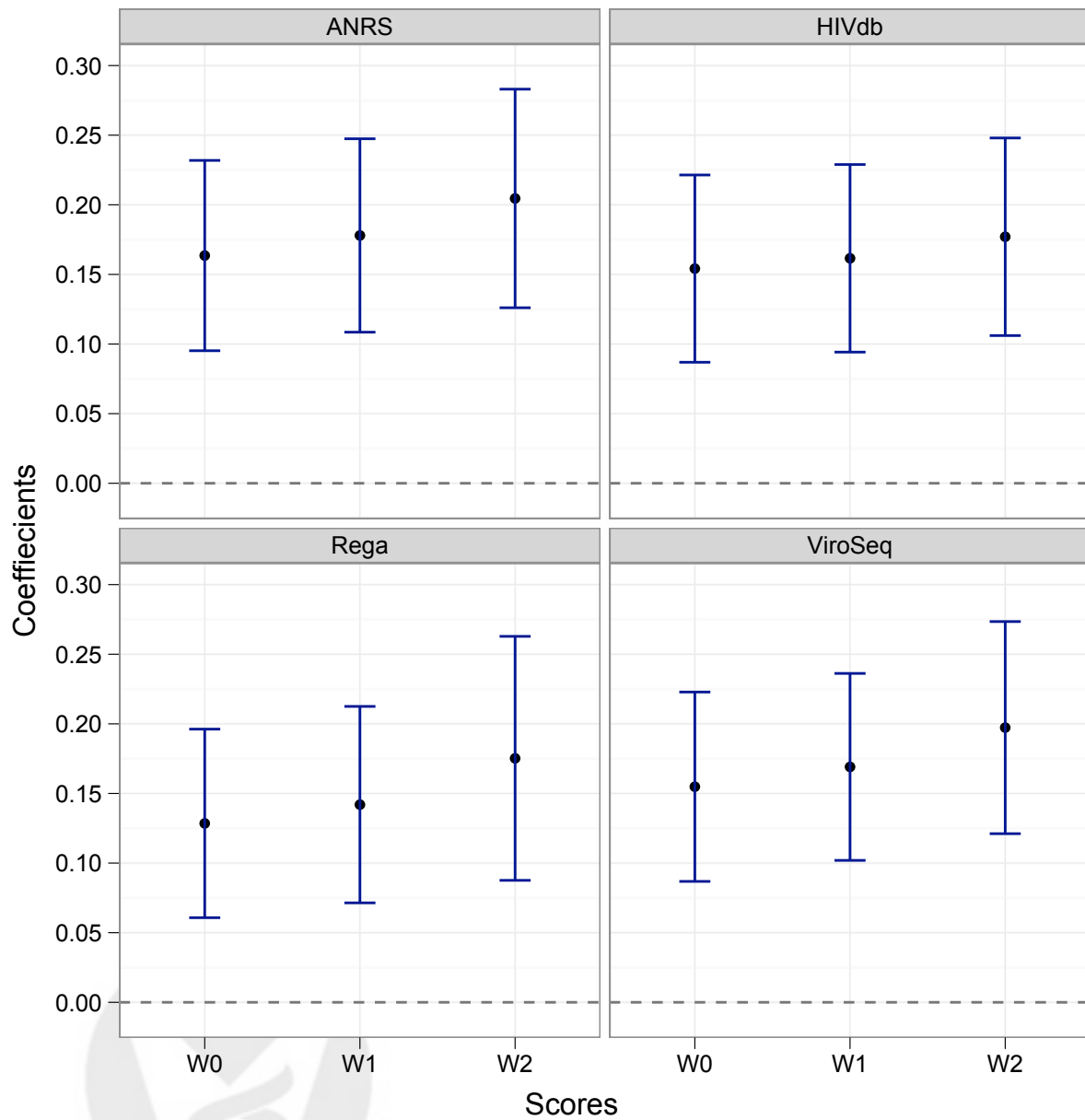
Figure 2: TMLE's of the coefficients for 4 algorithms, ANRS, HIVdb, Rega, and ViroSeq, under three different weighting schemes: "unweighted" (W0), "boosted PI weighted" (W1), and "comprehensive weighted" (W3). The brackets represent the 95% confidence intervals according the influence curve based estimate of the standard error.
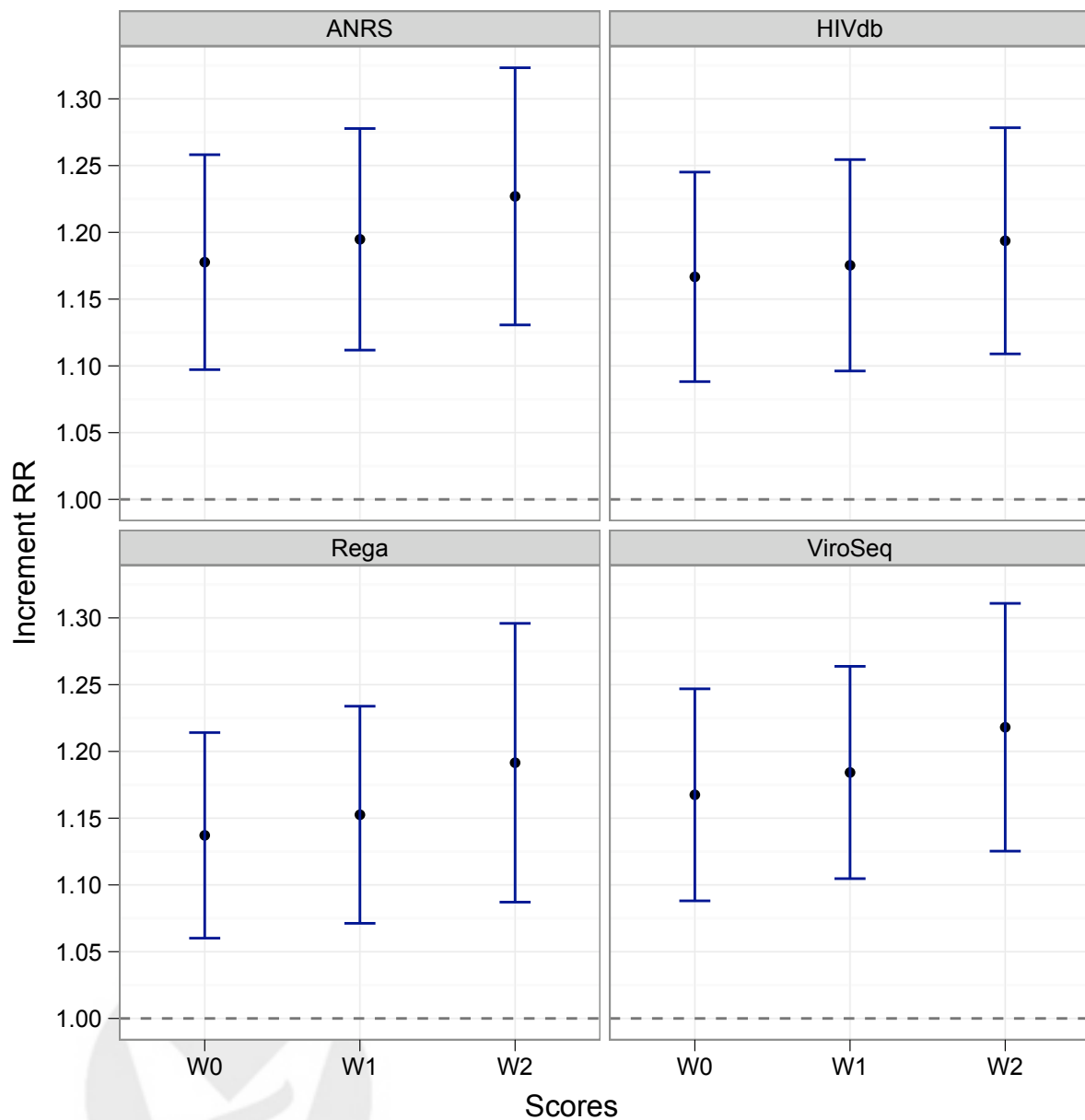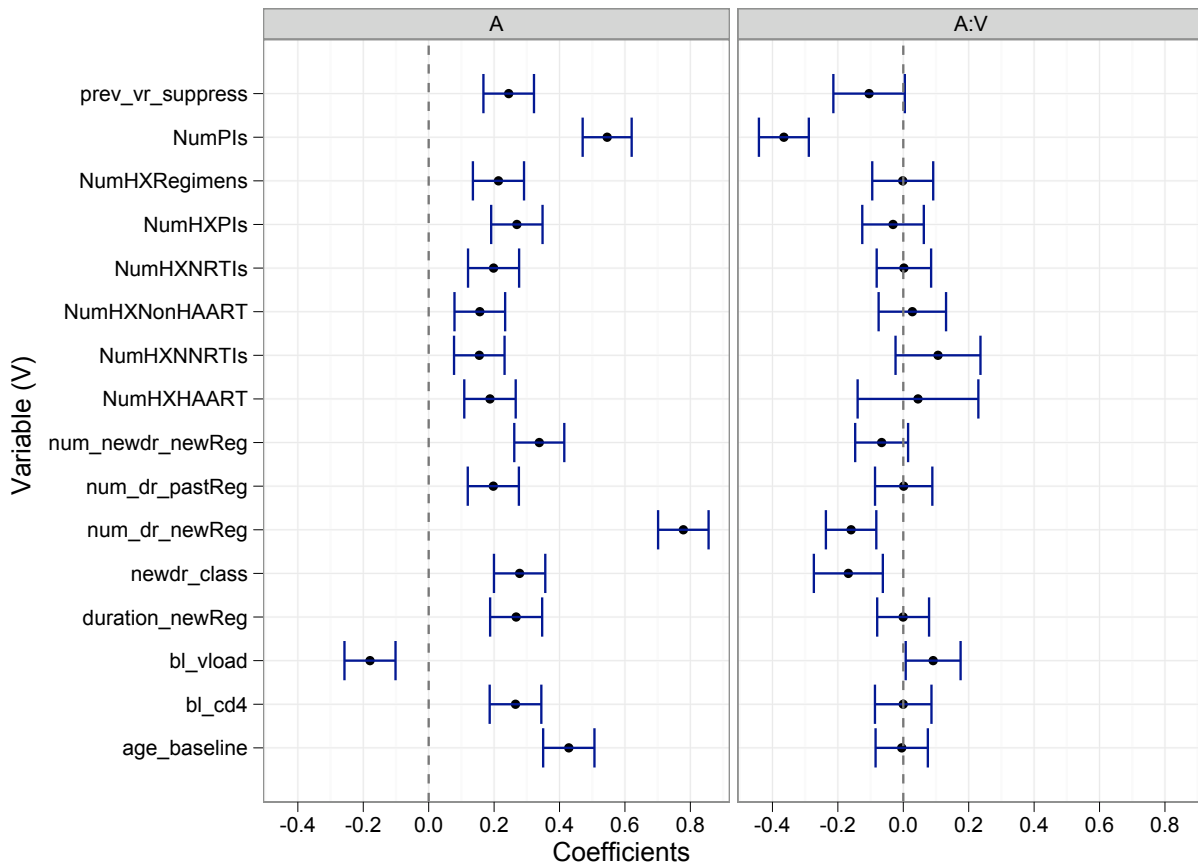
18

Figure 3: TMLE's of the increment RR for 4 algorithms, ANRS, HIVdb, Rega, and ViroSeq, under three different weighting schemes: "unweighted" (W0), "boosted PI weighted" (W1), and "comprehensive weighted" (W3). The brackets represent the 95% confidence intervals according the influence curve based estimate of the standard error using the delta method

19

Figure 4: TMLE's of the coefficients ($\beta$) for ANRS algorithm under comprehensive weighting adjusted by the other covariates. Left plot is the coefficient for the main effect, $A$, and the right plot is the coefficient of the effect modification or interaction term, $A\!:\!V$. The brackets represent the 95% confidence intervals according the influence curve based estimate of the standard error. The covariates adjusted are listed from top to bottom: history of virologic suppression prior to baseline (prev_vr_suppress), number of PIs in new regimen (NumPIs), number of regimens received prior to baseline (NumHXRegimens), number of PIs received prior to baseline (NumHXPIs), number of NRTIs received prior to baseline (NumHXNRTIs), number of non-HAART regimens received prior to baseline (NumHXNonHAART), number of NNRTIs received prior to baseline (NumHXNNRTIs), number of HAART regimens received prior to baseline (NumHXHAART), number of new ARVs in new regimen (num_newdr_newReg), number of ARVs received prior to baseline (num_dr_pastReg), number of ARVs in new regimen (num_dr_newReg), number of new ARV classes in new regimen (new_dr_class), duration of new regimen in weeks (duration_newReg), plasma HIV-1 RNA level at baseline in log-copies/ml (bl_vload), CD4 count at baseline in cells/ml (bl_cd4), and age at baseline in years (age_baseline)
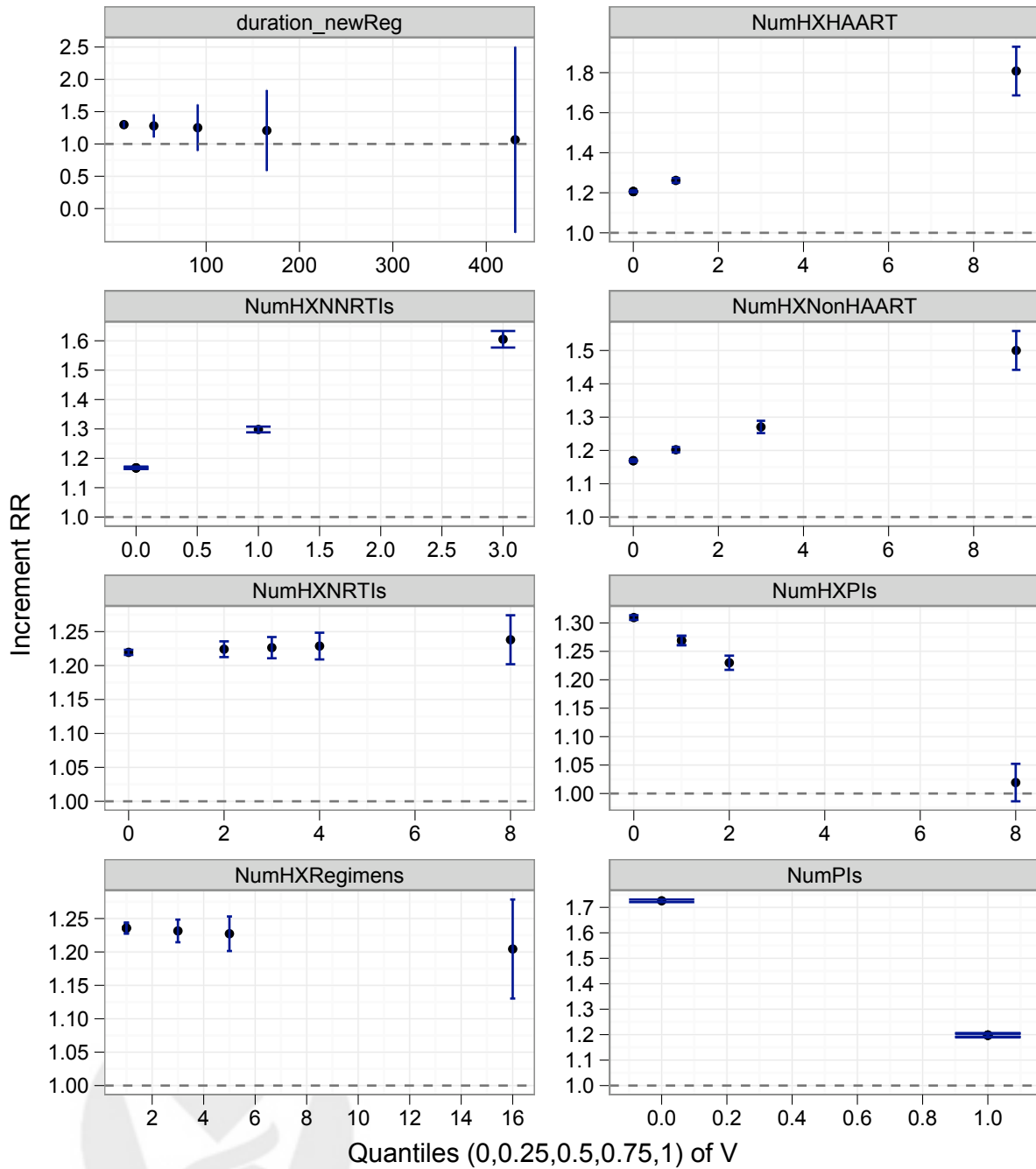
Figure 5: TMLE's of the increment RR for ANRS algorithm under comprehensive weighting modified by covariate $V$ at quantile levels $V_q = \{0\%, 25\%, 50\%, 75\%, 100\%\}$ for each covariate $V$. The brackets represent the 95% confidence intervals according the influence curve based estimate of the standard error. Covariate $V$ are as follows (left to right, top to bottom): duration of new regimen in weeks (duration_newReg), number of HAART regimens received prior to baseline (NumHXHAART), number of NNRTIs received prior to baseline (NumHXNNRTIs), number of non-HAART regimens received prior to baseline (NumHXNonHAART), number of NRTIs received prior to baseline (NumHXNRTIs), number of PIs received prior to baseline (NumHXPIs), number of regimens received prior to baseline (NumHXRegimens), and number of PIs in new regimen (NumPIs).
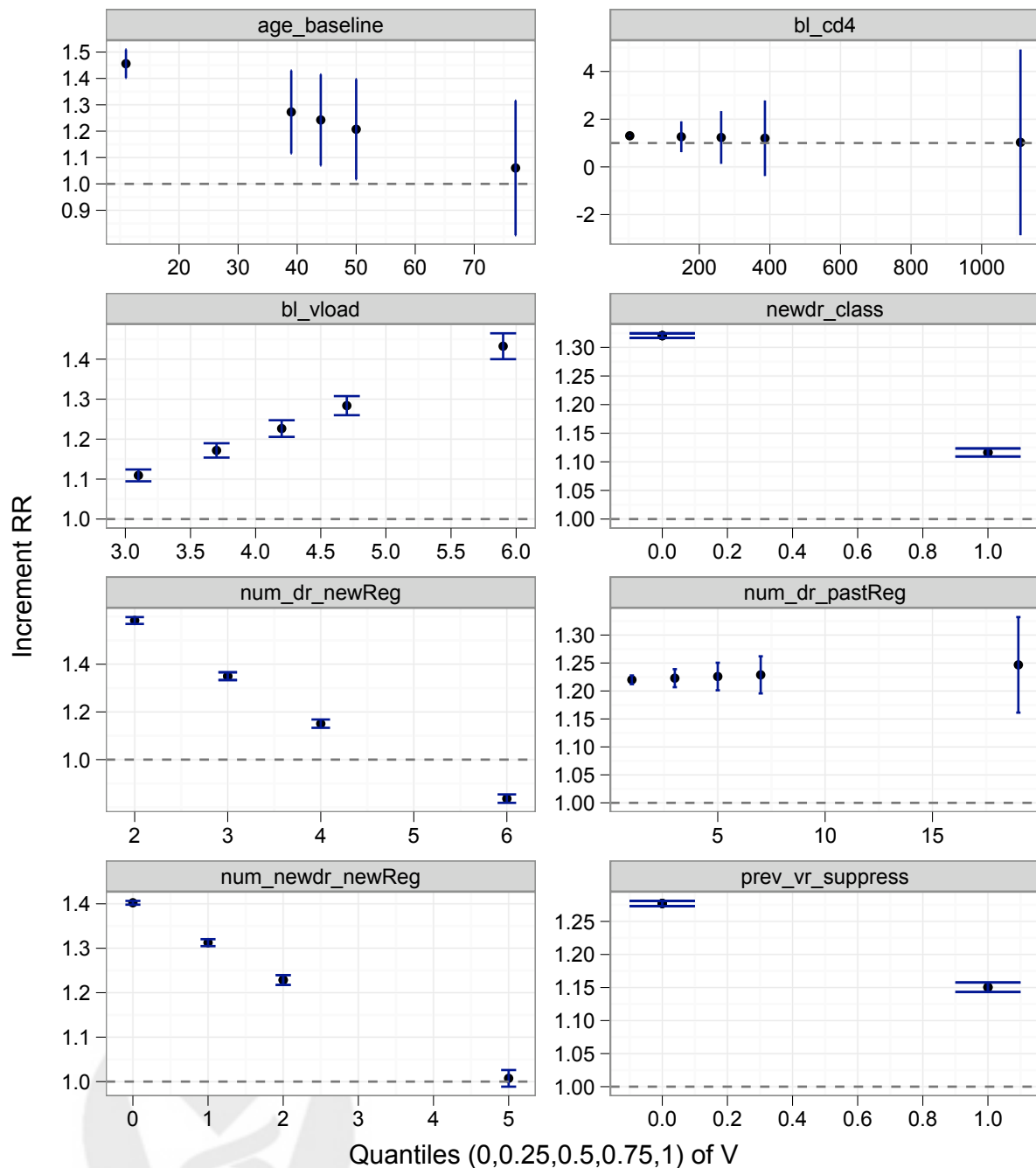
21

Figure 6: TMLE's of the increment RR for ANRS algorithm under comprehensive weighting modified by covariate $V$ at quantile levels $V_q = \{0\%, 25\%, 50\%, 75\%, 100\%\}$ for each covariate $V$. The brackets represent the 95% confidence intervals according the influence curve based estimate of the standard error. Covariate $V$ are as follows (left to right, top to bottom): age at baseline in years (age_baseline), CD4 count at baseline in cells/ml (bl_cd4), plasma HIV-1 RNA level at baseline in log-copies/ml (bl_vload), number of new ARV classes in new regimen (new_dr_class), number of ARVs in new regimen (num_dr_newReg), number of ARVs received prior to baseline (num_dr_pastReg), number of new ARVs in new regimen (num_newdr_newReg), history of virologic suppression prior to baseline (prev_vr_suppress).

22

# 9    Discussion

In this paper, we introduced three new estimators for the parameters of the conditional relative risk developed using targeted maximum likelihood methodology under a flexible multiplicative semi-parametric model. The most prevalent estimators in the literature, especially in the field of epidemiology, rely on parametric models. However, in a world where models are often incorrect, this reliance on full parametric models can result in biased estimates and inaccurate conclusions. TMLE's are developed to avoid reliance on these, often incorrect, models and reduce bias by targeting estimation towards the parameter of interest.

The TMLE's presented here are all double robust estimators for the parameter of interest, making them more resilient to the effects of model misspecification. The double robust property states that the resulting estimate is unbiased if either the estimator for $\bar{Q}_0(A, W)$ or the estimator for $g_0(A \mid W)$ is consistent. In the world clinical trials and controlled experiments, where $g_0(A \mid W)$ is known, this property is especially beneficial. However, coupled with data-adaptive methods, such as super learner (van der Laan et al., 2007), these TMLE's are also beneficial in the world of observational studies.

All three estimators are developed under a multiplicative semi-parametric model, which conveniently accommodates both binary and continuous covariates and effect modifiers, while adjusting for additional covariates flexibly using data-adaptive methods. This model, though more flexible than the fully parametric models commonly used in these studies, still requires that the analyst specify the components of the model relating to the variable of interest. In this paper, we discuss and implement two possible model forms, the main effect model and single effect modification model. It is important to note that though possible model forms are not restricted to these two models, they are restricted to models linear in the variable of interest ($A$). We also allow models with additional effect modifiers, but as the dimension of the fluctuation parameter increases, a sequential update may be required (Tuglus and van der Laan, 2010). In addition, the corresponding targeted likelihood also provides a framework for model selection among effect modifiers (Tuglus and van der Laan, 2011).

The three TMLE's introduced in this paper are defined and developed according to three different initial densities: log-binomial, Poisson, and overdispersed exponential. As mentioned previously, for parameters of the conditional relative risk, the log-binomial density is the natural choice. However, analogous to its parametric counterpart, the log-binomial-based TMLE is plagued by the instability inherent in log-binomial regression. As in previous studies before us (e.g. Zou (2004) and McNutt et al. (2003)), we turn to Poisson regression as a more reliable alternative. Although the Poisson-based TMLE does not correctly assume a binary outcome, it still possesses the double robust property, and its corresponding influence curve provides correct inference for the parameter of interest. We consider the Poisson-based TMLE the most practical TMLE of the three presented and therefore focus our efforts on its implementation and recommend it for general application.

In an effort to avoid reliance on distributional assumptions of $P_0(Y|A, W)$, we developed the TMLE based on the more general density from the overdispersed exponential family. This TMLE is developed directly from the overall efficient score for the semi-parametric multiplicative regression model and makes no assumptions on the distributional form of the residuals, $Y - E[Y|A, W]$. We show in Appendix A that the efficient scores of the log-binomial and Poisson-based TMLE can be derived directly from this more general overall efficient score. Although, this TMLE is a more general TMLE for the parameter of interest, it is not easily implemented using standard software packages, which makes this TMLE less practical for general implementation and application.

In this paper, we outline implementation instructions for the Poisson-based TMLE and present a simulation study in which we verify its double robust properties. We also demonstrate its performance under ETA violations and increases in confounding among the covariates. Similar to previously presented TMLE's developed under a semi-parametric model (Tuglus and van der Laan, 2008, 2010), this TMLE shows strong resilience to ETA violations. This is attributed to its lack of dependence on inverse probability weighting which is present in some of the previously developed TMLE's (van der Laan et al., September, 2009). Like its predecessors, this resilience does not nullify the positivity (or ETA) assumption from Section 4. The resilience comes from a natural extrapolation of the $P_0(A|W)$ mean process. When $A$ is continuous, the $P_0(A|W)$ mean process is estimated given the observed values of $A$ and $W$, and naturally extrapolates over

areas of lower support. These areas of lower support rely on the accuracy of the estimator for $P_0(A|W)$. Therefore in practice, it is important to acknowledge this reliance when strong ETA violations are present.

In application, we demonstrated the usefulness of the TMLE for estimation of the main effect of a particular variable, in this case a genotypic score, controlling for confounding This type of analysis is particularly useful in biomarker discovery and variable importance analyses when we want to accurately test the importance of many variables (see Tuglus and van der Laan (2008)). We also took it one step further and showed how the model can be augmented to test modifications of an association by a single or set of covariates, where the effect of a genotypic score on the relative risk of virologic response was modified by a variable such as the baseline viral load. Accurate and interpretable effect modification analysis such as this can be useful in clinical trials to test how the effect of a treatment ($A$) is modified by a particular gene expression, for instance.

In summary, conditional relative risk parameters are useful and interpretable in many epidemiology and medical studies. In this paper, we have shown that with a semi-parametric multiplicative model, one can obtain robust estimates of relative risk parameters, even under strong confounding and ETA violations. We introduced three possible TMLE's of parameters in the semi-parametric multiplicative model, and we implemented the more practical of these TMLE's - based on a Poisson density. We demonstrated this TMLE's DR properties, illustrated that it achieves proper inference and showed its improved performance over existing methods. We also illustrated the TMLE's value in a real data analysis. This TMLE can be applied using standard statistical software and will be useful in a wide variety of applications.

# References

Aluísio J D Barros and Vânia N Hirakata. Alternatives for logistic regression in cross-sectional studies: an empirical comparison of models that directly estimate the prevalence ratio. *BMC Med Res Methodol*, 3: 21, Oct 2003. doi: 10.1186/1471-2288-3-21.

A Buja, T Hastie, and R Tibshirani. Linear smoothers and additive models (with discussion). *Annals of Statistics*, 17(453-555), 1989.

RE Carter, SR Lipsitz, and BC Tilley. Quasi-likelihood estimation for relative risk regression models. *Biostatistics*, 6(1):39–44, 2005.

R Chen, W Hardle, O Linton, and E Severance-Lossin. Estimation and variable selection in additive non-parametric regression models. In W Hardle and M Schimek, editors, *Proceedings of the COMPSTAT Satellite Meeting Semmering 1994*, 1996.

Agence National de Recerche sur le SIDA (ANRS). Anrs genotypic resistence guidelines (verison 13). URL http://www.hivfrenchresistence.org/2007.

J Durant, P Clevenbergh, P Halfon, and et al. Drug-resistance genotyping in hiv-1 therapy: the viradapt randmised controlled trial. *Lancet*, 353:2195–99, 1999.

B. Efron, Hastie H., Johnstone, and Tibshirani. Least angle regression (with discussion). *Annals of Statistics*, 2003.

SH Eshleman, J Jr. Hackett, and P et al. Swanson. Performance of the celera diagnostics viroseq hiv-1 genotyping system for sequence-based analysis of diverse human immunodeficiency virus type 1 strains. *J Clin Microbiology*, 42(2711-7), 2004.

T. J. Hastie and R.J. Tibshirani. *Generalized Additive Models*. Chapman and Hall, London, 1990.

J. Lee. Odds ratio or relative risk for cross-sectional data? *Int J Epidemiol*, 23:201–3, 1994.

J Lee and KS Chia. Estimation of prevalence rate ratios for cross-sectional data: an example in occupational epidemiology. *Br J Ind Med*, 50:861–2, 1993.

TF Liu and RW. Shafer. Web resources for hiv type 1 genotypic-resistence test interpretation. *Clin Infect Dis*, 42:1608–18, 2006.

AR Localio, DJ Margolis, and JA Berlin. Relative risks and confidence intervals were easily computed indirectly from multivariate logisitic regression. *Journal of Clinical Epidemiology*, 60:874–882, 2007.

Thomas Lumley, Richard Kronmal, and Shuangge Ma. Relative risk regresssion in medical research: Models, contrasts, estimators, and algorithms. Technical Report 293, University of Washington, 2006.

P. McCullagh and J. A. Nelder. *Generalized Linear Models*. Number 37 in Monographs on Statistics and Applied Probability. Chapman & Hall, 2nd edition, 1989.

Louise-Anne McNutt, Chuntao hu, Xiaonana Xue, and Paul Hafner. Estimating the relative risk in cohort studies and clinical trials of common outcomes. *Am J Epidemiol*, 157(10):940–943, 2003.

R. Neugebauer and M.J. van der Laan. Why prefer double robust estimates. *Journal of Statistical Planning and Inference*, 129(1-2):405–26, 2005.

US Department of Health and Human Services Panel on Clinical Practices for Treatment of HIV Infection A. Guidelines for the use of antiretroviral agents in hiv-1-infected adults and adolescents (the living document, june 26, 2010). URL http://www.aidsinfo.nih.gov/guidelines/.

SY Rhee, WJ Fessel, TF Liu, NM Marlowe, CM Rowland, RA Rode, AM Vandamme, K Van Laethem, F Brun-Vezinet, V Calvez, J Taylor, L Hurley, M Horberg, and RW Shafer. Predictive value of hiv-1 genotypic resistance test interpretation algorithms. *J Infect. Dis.*, 200(3):453–63, 2009.

J.M. Robins. Robust estimation in sequentially ignorable missing data and causal inference models. In *Proceedings of the American Statistical Association: Section on Bayesian Statistical Science*, pages 6–10, 1999.

J.M. Robins. A new approach to causal inference in mortality studies with sustained exposure periods - application to control of the healthy worker survivor effect. *Mathematical Modelling*, 7:1393–1512, 1986.

J.M. Robins. Addendum to: "A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect" [Math. Modelling **7** (1986), no. 9-12, 1393–1512; MR 87m:92078]. *Comput. Math. Appl.*, 14(9-12):923–945, 1987. ISSN 0097-4943.

*Estimation of prevalence ratios when PROC GENMOD does not converge*, number 270-28, Cary, NC, 2003. SAS Institute, Proceedings of the 28th Annual SAS Users Group International Conference, March 30-April 2, 2003.

T Severini and W Wong. Profile likelihood and conditionally parametric models. *Annals of statistics*, 20:1768–1802, 1992.

TA Severini and JG Staniswalis. Quasi-likelihood estimation in semiparametric models. *Journal of the American Statisical Association*, 89:501–511, 1994.

S. Sinisi and M.J. van der Laan. The deletion/substitution/addition algorithm in loss function based estimation: Applications in genomics. *Journal of Statistical Methods in Molecular Biology*, 3(1), 2004.

T Skov, J.A. Deddens, and M.R. Petersen. Prevalence proportion ratios: estimation and hypothesis testing. *Int J Epidemiol*, 27:91–95, 1998.

PE Speckman. Regression analysis for partially linear models. *Journal of the Royal Statistical Society, Series B*, 50:413–436, 1988.

T Stijnen and HC Van Houwelingen. Relative risk, risk difference and rate difference models for sparse stratified data: A pseudo likelihood approach. *Statistics in Medicine*, 12(24):2285–2303, 1993.

R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2010.

C. Tuglus and M.J. van der Laan. Targeted methods for biomarker discovery, the search for a standard. *UC Berkeley Working Paper Series*, page http://www.bepress.com/ucbbiostat/paper233/, 2008.

C. Tuglus and M.J. van der Laan. Repeated measures semiparametric regression using targeted maximum likelihood methodology with application to transcription factor activity discovery. *UC Berkeley Working Paper Series*, page http://www.bepress.com/ucbbiostat/paper261, 2010.

C. Tuglus and M.J. van der Laan. *Robust Semiparametric Regression Estimation Using Targeted Maximum Likelihood with Application to Biomarker Discovery and Epidemiology*. PhD thesis, UC Berkeley, 2011.

C Tural, L Ruiz, C Holtzer, and et al. Clinical utility of hiv-1 genotyping and expert advice: the havana trial. *AIDS*, 16:209–18, 2002.

M.J. van der Laan. Statistical inference for variable importance. *International Journal of Biostatistics*, 2 (1), 2006.

M.J. van der Laan and D. Rubin. Targeted maximum likelihood learning. *The International Journal of Biostatistics*, 2(1), 2006.

M.J. van der Laan, E. Polley, and A. Hubbard. Super learner. *Statistical Applications in Genetics and Molecular Biology*, 6(25), 2007. ISSN 1.

M.J. van der Laan, S. Rose, and S. Gruber. Readings on targeted maximum likelihood estimation. *Technical report, working paper series http://www.bepress.com/ucbbiostat/paper254*, September, 2009.

K Van Laethem, A De Luca, A Antinori, A Cingolani, CF Perna, and AM Vandamme. A genotypic drug resistence interpretation algorithm that significantly predicts therapy response in hiv-1 infected patients. *Antivir Ther*, 7:123–9, 2002.

S Wacholder. Binomial regression in glim, esitmating risk ratios and risk differences. *Am J Epidemiol*, 123: 174–84, 1986.

C Zocchett, Dario Consonni, and Pier Alberto Bertazzi. Re: Estimation of prevalence rate ratios from cross-sectional data (letter). *Int J Epidemiol*, 24:1064–1105, 1995.

G. Zou. A modified poisson regresion approach to prospective studies with binary data. *Am J Epidemiol*, 159:702–6, 2004.

# A    Efficient Influence curve for effect parameter of multiplicative semi-parametric model

We present the efficient influence curve for the multiplicative semi-parametric model parameter, and show that it can be represented as a score of the fluctuation in the overdispersed exponential family of canonical form. We first present below the general form of the efficient influence curve, then present the overdispersed exponential family of canonical form, and finally construct a submodel in this exponential family with score equal to the efficient influence curve.

26

## A.1 Efficient Influence curve for effect parameter of multiplicative semi-parametric model

Assume a multiplicative semi-parametric model of the following form
$\bar{Q}(A,W)m_\beta^\star(A,V) = \theta(W)$ under the following constraints: $m^\star(A = 0, V|\beta) = 1$ and
$0 \le m_\beta^\star(A,V):$ for all $\{a,v\} \in \{A,V\}$, where $m_\beta^\star(A,V) = e^{-m_\beta(A,V)}$. The effect parameter of interest is defined as $\Psi(P) = \beta$, where $\Psi(P_0) = \beta_0$ is the true parameter defined under the true data generating distribution.

As presented in van der Laan (2006), the orthogonal complement of the nuisance tangent space for estimation of $\beta$ is found to be of the form.

$$T_{nuis}^\perp(P_0) = \{h_{\bar{Q}_0,g_0}(A,W) - E_0[h_{\bar{Q}_0,g_0}(A,W)|W]\}(Ym_{\beta_0}^\star(A,V) - \theta_0(W))$$

with class of estimating functions

$$(O, \beta_0, \bar{Q}_0, g_0) \to D_{h,\bar{Q}_0,g_0}(O|\beta_0) \equiv \{h_{\bar{Q}_0,g_0}(A,W) - E_0[h_{\bar{Q}_0,g_0}(A,W)|W]\}(Ym_{\beta_0}^\star(A,V) - \theta_0(W))$$

for $\beta_0$ indexed by $h$, where $\theta_0 = E_0(Y|A = 0, W)$, and $g_0 = P_0(A|W)$. The corresponding influence curve is defined as

$$IC_{h,\bar{Q}_0,g_0}(O) = -\frac{D_{h,\bar{Q}_0,g_0}(O|\beta_0)}{\frac{d}{d\beta_0}E_0[D_{h,\bar{Q}_0,g_0}(O|\beta_0)]}.$$

The optimal choice of $h$, $h_{opt}$, is such that for any vector c,

$$c^T Cov(IC_{h_{opt},\bar{Q}_0,g_0})c \le c^T Cov(IC_{h,\bar{Q}_0,g_0})c$$

for all possible $h_{\bar{Q}_0,g_0}(A,W)$, thus providing the most efficient estimating function. For effect parameter $\beta$ under the presented multiplicative semi-parametric model, $h_{opt,\bar{Q}_0,g_0}$ is defined as follows.

Define the following terms.
$$H_0(O|\beta_0) \equiv Ym_{\beta_0}^\star(A,V)$$

$$\epsilon(\beta_0) \equiv H_0(O|\beta_0) - E_0[H_0(O|\beta_0)|W]$$

$$\epsilon'(\beta_0|A,W) \equiv \frac{d}{d\beta}E_0[\epsilon(\beta)|A,W]\bigg|_{\beta=\beta_0}$$

$$\sigma^2(A,W) \equiv E_0(\epsilon^2(\beta_0)|A,W)$$

where,

$$h_{opt,\bar{Q}_0,g_0}(A,W) = \frac{1}{\sigma^2(A,W)}\left\{\epsilon'(\beta_0|A,W) - \frac{\int \frac{\epsilon'(\beta_0|A,W)}{\sigma^2(A,W)}dP_0(a|W)}{\int \frac{1}{\sigma^2(A,W)}dP_0(a|W)}\right\}$$

or

$$h_{opt,\bar{Q}_0,g_0}(A,W) = \frac{1}{\sigma^2(A,W)}\left\{\epsilon'(\beta_0|A,W) - \frac{E_0\left[\frac{\epsilon'(\beta_0|A,W)}{\sigma^2(A,W)}\Big|W\right]}{E_0\left[\frac{1}{\sigma^2(A,W)}\Big|W\right]}\right\}.$$

This can be rewritten as

$$h_{opt,\bar{Q}_0,g_0}(A,W) = \frac{1}{\sigma^2(A,W)}\left\{\bar{Q}_0(A,W)\frac{d}{d\beta}m_{\beta_0}^\star(A,V) - \frac{E_0\left[\frac{\bar{Q}_0(A,W)\frac{d}{d\beta}m_{\beta_0}^\star(A,V)}{\sigma^2(A,W)}\Big|W\right]}{E_0\left[\frac{1}{\sigma^2(A,W)}\Big|W\right]}\right\}$$

This results in the following efficient score/efficient estimating function:

$$D_{h_{opt},\bar{Q}_0,g_0}(A,W|\beta_0) = h_{opt,\bar{Q}_0,g_0}(A,W)(Ym_{\beta_0}^\star(A,V) - \theta_0(W)).$$

27

Note that $\sigma^2(A,W) = m_{\beta_0}^\star(A,V)^2 \sigma_Y^2(A,W)$, where $\sigma_Y^2(A,W) = Var(Y|A,W)$, so that the efficient estimating function can be simplified further as

$$D_{h,\bar{Q}_0,g_0}(A,W|\beta_0) = h_{opt,\bar{Q}_0,g_0}^*(A,W)(Y - \bar{Q}_0(A,W))$$

$$h_{opt,\bar{Q}_0,g_0}^*(A,W) = \frac{\bar{Q}_0(A,W)}{\sigma_Y^2(A,W)}\left\{\frac{d}{d\beta}m_{\beta_0}(A,V) - \frac{E_0\left[\frac{d}{d\beta}m_{\beta_0}(A,V)\frac{1}{\sigma^2(A,W)}\Big|W\right]}{E_0\left[\frac{1}{\sigma^2(A,W)}\Big|W\right]}\right\},$$

which can be rewritten as

$$D_{h_{opt,\bar{Q}_0,g_0},\bar{Q}_0,g_0}(A,W|\beta_0) = h_{opt}^*(A,W)(Y - \bar{Q}_0(A,W)) \tag{4}$$

$$h_{opt,\bar{Q}_0,g_0}^*(A,W) = \frac{\bar{Q}_0(A,W)}{\sigma_Y^2(A,W)}\left\{\frac{d}{d\beta}m_{\beta_0}(A,V) - \frac{E_0\left[\frac{d}{d\beta}m_{\beta_0}(A,V)\frac{\bar{Q}_0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E_0\left[\frac{\bar{Q}_0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}\right\}. \tag{5}$$

Up until this point no distributional assumptions are made about the conditional distribution of $Y$. The above efficient estimating function can be simplified further by assuming a form for $\sigma_Y^2(A,W) = Var(Y|A,W)$. If one assumes a bernoulli or binomial outcome, $\sigma_Y^2(A,W) = \bar{Q}_0(A,W)(1 - \bar{Q}_0(A,W))$, and the above estimating function reduces to the estimating function used in Section 5.1.1 for the targeted maximum likelihood update of an initial log-binomial density. If one assumes a Poisson count outcome, then $\sigma_Y^2(A,W) = \bar{Q}_0(A,W)$, and the above estimating function reduces to the estimating function used in Section 5.1.2 for the targeted maximum likelihood update of an initial Poisson density.

## A.2  Overdispersed exponential family

The density of the overdispersed exponential family in canonical form is represented as follows

$$P_{\eta,\tau}^0(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\eta Y - B(\eta)}{d(\tau)}\right\}.$$

We define $\eta = \bar{Q}^0(A,W) = \theta(W)\exp(m_{\beta^0}(A,V))$ and define a class of submodels, fluctuated by parameter $\epsilon$ as

$$P_{\eta,\tau}^0(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\eta(\epsilon)Y - B(\eta(\epsilon))}{d(\tau)}\right\},$$

where $\eta(\epsilon) = \bar{Q}(\epsilon)(A,W) = \theta(\epsilon)(W)\exp(m_{\beta(\epsilon)}(A,V))$, $\theta(\epsilon)(W) = \theta(W)\exp(\epsilon r_{\bar{Q}^0,g^0}(W))$ and $\beta(\epsilon) = \beta + \epsilon$. Therefore the score of the above likelihood with respect to epsilon is as follows

$$\frac{d}{d\epsilon}\log P_{\eta,\tau}^0(\epsilon)(Y|A,W) = \frac{1}{d(\tau)}\left\{\frac{d}{d\epsilon}\eta(\epsilon)Y - B'(\eta(\epsilon))\frac{d}{d\epsilon}\eta(\epsilon)\right\},$$

which can be rearranged as

$$\frac{d}{d\epsilon}\log P_{\eta,\tau}^0(\epsilon)(Y|A,W) = \frac{1}{d(\tau)}\frac{d}{d\epsilon}\eta(\epsilon)(Y - B'(\eta(\epsilon))).$$

By definition of the canonical form, $B'(\eta(\epsilon)) = \bar{Q}(\epsilon)(A,W)$, so that the score can be written as

$$\frac{d}{d\epsilon}\log P_\tau^0(\epsilon)(Y|A,W) = \frac{1}{d(\tau)}\frac{d}{d\epsilon}\bar{Q}(\epsilon)(A,W)(Y - \bar{Q}(\epsilon)(A,W)).$$

Given the form for $\bar{Q}(\epsilon)(A,W)$, it follows that

$$\frac{d}{d\epsilon}\bar{Q}(\epsilon)(A,W) = \bar{Q}(\epsilon)(A,W)\left\{\frac{d}{d\beta}m_\beta(A,V) + r_{\bar{Q}^0,g^0}(W)\right\},$$

28

so that the score at $\epsilon = 0$ is given by

$$\frac{d}{d\epsilon} \log P_\tau^0(\epsilon)(Y|A,W)\Big|_{\epsilon=0} = \frac{\bar{Q}^0(A,W)}{d(\tau)} \left\{ \frac{d}{d\beta} m_\beta(A,V) + r_{\bar{Q}^0,g^0}(W) \right\} (Y - \bar{Q}^0(A,W)).$$

This score is equivalent to the form shown in Equation (6) when $d(\tau) = \sigma_Y^2(A,W)$. This proves that the estimating function/efficient score presented above is a score of a likelihood in the overdispersed exponential family. This provides another verification of the fact that the claimed efficient score is indeed an efficient score.

# B  Derivation of TMLE's

## B.1  Log-binomial

Under the Log-binomial distribution, the initial density is a binomial density defined as

$$P^0(Y|A,W) = \bar{Q}^0(A,W)^Y (1 - \bar{Q}^0(A,W))^{1-Y},$$

where $\bar{Q}^0(A,W) = P^0(Y = 1|A,W) = \theta^0(W)e^{m_{\beta^0}(A,V)}$ with the associated fluctuation

$$P^0(\epsilon)(Y|A,W) = \bar{Q}^0(\epsilon)(A,W)^Y (1 - \bar{Q}^0(\epsilon)(A,W))^{1-Y},$$

given $\bar{Q}^0(\epsilon)(A,W) = \theta^0(\epsilon)(W)e^{m_{\beta^0(\epsilon)}(A,V)}$

The associated score for the above likelihood with respect to $\epsilon$ at $\epsilon = 0$ is as follows

$$S(r) = \frac{1}{1 - \bar{Q}^0(A,W)} \left\{ \frac{d}{d\beta} m_{\beta^0}(A,V) + r_{\bar{Q}^0,g^0}(W) \right\} (Y - \bar{Q}^0(A,W)).$$

The efficient score for effect parameter $\beta_0$ under the multiplicative semi-parametric model of Section 3 assuming a log-binomial distribution is defined below under $P_0$

$$D_{h_0,\bar{Q}_0}(O) = h_{\bar{Q}_0,g_0}(A,W)(Y - \bar{Q}_0(A,W)),$$

where

$$h_{\bar{Q}_0,g_0}(A,W) = \frac{1}{1 - \bar{Q}_0(A,W)} \left\{ \frac{d}{d\beta_0} m_{\beta_0}(A,V) - \frac{E_0\left[ \frac{\bar{Q}_0(A,W)}{1-\bar{Q}_0(A,W)} \frac{d}{d\beta_0} m_{\beta_0}(A,V) \Big| W \right]}{E_0\left[ \frac{\bar{Q}_0(A,W)}{1-\bar{Q}_0(A,W)} \Big| W \right]} \right\}.$$

This score is shown to belong to the general class of estimating functions for effect parameter $\beta$ under the multiplicative semi-parametric model in Appendix A, equals the efficient score for a binary outcome.

The proper form of the fluctuation function $r_{\bar{Q}^0,g^0}(W)$ is therefore as follows

$$r_{\bar{Q}^0,g^0}(W) = -\frac{E\left[ \frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)} \frac{d}{d\beta^0} m_{\beta^0}(A,V) \Big| W \right]}{E\left[ \frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)} \Big| W \right]}.$$

Given a model $m_\beta(A,V)$ that is linear in $\beta$, the model,

$$\log \bar{Q}^0(\epsilon)(A,W) = m_{\beta^0(\epsilon)}(A,V) + \log \theta^0(\epsilon)(W),$$

can be rearranged as an update to the initial fit

$$\log \bar{Q}^0(\epsilon)(A,W) = \log \bar{Q}^0(A,W) + \epsilon^T \frac{d}{d\beta^0} m_{\beta^0}(A,V) + \epsilon^T r_{\bar{Q}^0,g^0}(W).$$

29

Therefore the update can be achieved by estimating $\epsilon$ with standard maximum likelihood estimation. The update can be completed using log-binomial regression setting the initial estimate, $\bar{Q}^0(A,W)$, as an offset and regressing $Y$ onto the following "clever covariate",

$$H^*_{\bar{Q}^0,g^0}(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\frac{d}{d\beta^0}m_{\beta^0}(A,V)\Big|W\right]}{E\left[\frac{\bar{Q}^0(A,W)}{1-\bar{Q}^0(A,W)}\Big|W\right]}.$$

## B.2 Poisson

Under the Poisson distribution, the initial density is a Poisson density defined as

$$P^0(Y|A,W) = \frac{\bar{Q}^0(A,W)^Y}{Y!}e^{-\bar{Q}^0(A,W)},$$

with the associated fluctuation

$$P^0(\epsilon)(Y|A,W) = \frac{\bar{Q}^0(\epsilon)(A,W)^Y}{Y!}e^{-\bar{Q}^0(\epsilon)(A,W)}.$$

The associated score for the above likelihood with respect to $\epsilon$ at $\epsilon = 0$ is as follows

$$S(r) = \left\{\frac{d}{d\beta}m_{\beta^0}(A,V) + r_{\bar{Q}^0,g^0}(W)\right\}(Y - \bar{Q}^0(A,W)).$$

The efficient score associated under the multiplicative semi-parametric model of Section 3 assuming a Poisson distribution is defined below under $P_0$. This score is also shown to belong to the general class of estimating functions for the effect parameter $\beta$ under the multiplicative semi-parametric model in Appendix A, and is the efficient score for this effect parameter $\beta_0$ given a Poisson (i.e. count) outcome. We have

$$D_{h_0,\bar{Q}_0}(O) = h_{\bar{Q}_0,g_0}(A,W)(Y - \bar{Q}_0(A,W)),$$

where

$$h_{\bar{Q}_0,g_0}(A,W) = \frac{d}{d\beta_0}m_{\beta_0}(A,V) - \frac{E_0[\frac{d}{d\beta_0}m_{\beta_0}(A,V)e^{m_{\beta_0}(A,V)}|W]}{E_0[e^{m_{\beta_0}(A,V)}|W]}.$$

It follows that the proper form for $r_{\bar{Q}^0,g^0}(W)$ is

$$r_{\bar{Q}^0,g^0}(W) = -\left\{\frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E[\frac{d}{d\beta^0}m_{\beta^0}(A,V)e^{m_{\beta^0}(A,V)}|W]}{E[e^{m_{\beta^0}(A,V)}|W]}\right\}.$$

Given a simple linear model for $m_\beta(A,V) = \beta A$ the above can rewritten as

$$r_{\bar{Q}^0,g^0}(W) = -\left\{A - \frac{E[Ae^{\beta^0 A}|W]}{E[e^{\beta^0 A}|W]}\right\}.$$

Similar to the log-binomial case, the update can be achieved by estimating $\epsilon$ with standard maximum likelihood estimation. The update is completed using Poisson regression with an offset equal to the initial fit and "clever covariate" defined as

$$H^*_{\bar{Q}^0,g^0}(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \left\{\frac{E[\frac{d}{d\beta^0}m_{\beta^0}(A,V)e^{m_{\beta^0}(A,V)}|W]}{E[e^{m_{\beta^0}(A,V)}|W]}\right\}.$$

30

## B.3  General semi-parametric multiplicative model

Assuming only the semi-parametric multiplicative model, we define the initial density as a member of the overdispersed exponential family as outlined in Appendix A such that

$$P_\tau^0(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\bar{Q}^0 Y - B(\bar{Q}^0)}{d(\tau)}\right\}.$$

and we define a class of submodels, fluctuated by parameter $\epsilon$ as

$$P_\tau^0(\epsilon)(Y|A,W) = h_c(Y,\tau)\exp\left\{\frac{\bar{Q}^0(\epsilon)Y - B(\bar{Q}^0(\epsilon))}{d(\tau)}\right\},$$

where $\bar{Q}(\epsilon)(A,W) = \theta(\epsilon)(W)\exp(m_{\beta(\epsilon)}(A,V))$, $\theta(\epsilon)(W) = \theta(W)\exp(\epsilon r_{\bar{Q}^0,g^0}(W))$ and $\beta(\epsilon) = \beta + \epsilon$. The score of the above likelihood with respect to $\epsilon$ defined at $\epsilon = 0$ is as follows

$$S_\tau(r) = \frac{\bar{Q}^0(A,W)}{d(\tau)}\left\{\frac{d}{d\beta}m_\beta(A,V) + r_{\bar{Q}^0,g^0}(W)\right\}(Y - \bar{Q}^0(A,W)).$$

The efficient score for the effect parameter $\beta_0$ only assuming the multiplicative semi-parametric model defined in Section 3 as shown previously in Appendix A is restated here

$$D_{h_{opt,\bar{Q}_0,g_0},\bar{Q}_0,g_0}(A,W|\beta_0) = h_{opt}^*(A,W)(Y - \bar{Q}_0(A,W)) \tag{6}$$

$$h_{opt,\bar{Q}_0,g_0}^*(A,W) = \frac{\bar{Q}_0(A,W)}{\sigma_Y^2(A,W)}\left\{\frac{d}{d\beta}m_{\beta_0}(A,V) - \frac{E_0\left[\frac{d}{d\beta}m_{\beta_0}(A,V)\frac{\bar{Q}_0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E_0\left[\frac{\bar{Q}_0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}\right\}. \tag{7}$$

It follows directly that

$$r_{\bar{Q}^0,g^0}(W) = -\left\{\frac{E\left[\frac{d}{d\beta^0}m_{\beta^0}(A,V)\frac{\bar{Q}^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E\left[\frac{\bar{Q}^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}\right\},$$

given $d(\tau) = \sigma_Y^2(A,W)$. Therefore the "clever covariate" can be defined as

$$H_{\bar{Q}^0,g^0}^*(A,W) = \frac{d}{d\beta^0}m_{\beta^0}(A,V) - \frac{E\left[\frac{d}{d\beta^0}m_{\beta^0}(A,V)\frac{\bar{Q}^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}{E\left[\frac{\bar{Q}^0(A,W)^2}{\sigma_Y^2(A,W)}\Big|W\right]}.$$

31