

Analysis of randomized trials with non-compliance and a binary outcome

James Robins¹ and Andrea Rotnitzky²

Harvard School of Public Health, Departments of ¹Epidemiology and ^{1,2}Biostatistics
and ²Di Tella University, Department of Economics

655 Huntington Ave, Boston, USA

¹robins@hsph.harvard.edu and ²andrea@hsph.harvard.edu

In this presentation we summarize results in Robins and Rotnitzky (2001) concerning the problem of estimating the actual effect of treatment in randomized clinical trials with noncompliance and a dichotomous outcome. We allow both for the assigned and actual treatment to be continuous categorical or ordinal and for the possibility that the assigned treatment has a direct effect on the outcome through pathways other than actual treatment. Our proposed methodology is based on the logistic structural mean model of Robins (1997) and Robins et. al. (1999).

Special challenges arise for analyzing binary outcomes with logistic structural mean models. Specifically, for continuous or count outcomes, unbiased estimating functions for the (structural) treatment effect of the additive and multiplicative structural mean models of Robins (1994) exist even in the presence of many continuous and discrete exogenous baseline covariates, provided the randomization probabilities (i.e. the conditional density of the instrument given the exogenous baseline covariates) are known. However, unless the causal sharp null hypothesis of no effect of either assigned or actual treatment holds, such unbiased estimating functions do not exist for the structural parameters of the logistic structural mean model. As a consequence, the best that can be hoped for are estimators, such as those proposed in this presentation, that are guaranteed both (i) to consistently estimate the (null) treatment effect when the null hypothesis of no effect of either assigned or actual treatment is true and, (ii) to have small bias when the true treatment effect is close but not equal to zero. In particular, one can use an estimator satisfying (i) in conjunction with its estimated standard error to construct a test that is an asymptotically distribution free α level test both of the intention to treat null hypothesis that the outcome is conditionally independent of assigned treatment given baseline covariates and of the null hypothesis of no causal effect of either assigned or actual treatment. In the presence of noncompliance the possibility of constructing estimators that satisfy conditions (i) and (ii) is one of the most important benefits to be realized from a randomized experiment. Because of the considerable cost associated with conducting a randomized experiment, we would regard as less than satisfactory any analytic approach that failed to satisfy (i) and (ii). Our work was stimulated by a recent paper on placebo controlled trials by Vansteelandt and Goetghebeur (2001) in which these authors, in conjunction with J Robins, proposed an estimation method that satisfied (i) and (ii). However, when some subjects in the placebo arm take the active treatment Vansteelandt and Goetghebeur modelling strategy can be internally inconsistent. For this reason, we prefer our proposed approach to that of Vansteelandt and Goetghebeur, in spite of their approach being less computationally demanding.

Under the exclusion assumption that assigned treatment has no direct effect on the outcome, the problem treated in this paper is equivalent to the problem of how to use an instrumental variable (i.e. assigned treatment) to estimate the effect of an endogenous (i.e. actual) treatment on a dichotomous response variable. This latter problem is a long standing problem in econometrics that was recently reviewed by Angrist (2001). When data on many continuous and discrete ex-

ogenous baseline covariates have been obtained, to the best of our knowledge, the only proposed econometric solutions are for the special case in which both assigned and actual treatment are themselves dichotomous, there are no defiers (i.e. there are no subjects such that they take treatment if and only if they are assigned to the control arm), the known randomization probabilities may depend on baseline covariates, and interest is in the treatment effect among those who comply with their assigned treatment (Abadie, 2001; Hirano, Imbens, Rubin and Zhou, 2000). For this special case, Abadie's estimator, in contrast to Hirano et al's, satisfies conditions (i) and (ii). In essence, the reason for the failure of Hirano's et al estimator is that it does not use knowledge of the randomization probabilities and thus, as shown by Robins and Ritov (1997), it cannot satisfy (i).

Often it is logistically or ethically impossible to conduct a randomized experiment. In that case one can sometimes find a variable that one is willing to consider an instrumental variable within levels of a vector of covariates. For example, one might regard the prescribed dose of a drug therapy to be an instrument for the dose actually taken within levels of the often high dimensional vector of prognostic factors used by physicians to determine the need for the therapy. Unlike in designed experiments, it would be typical for the probability of the instrument to depend in a complex fashion on such a high dimensional vector of covariates. Furthermore, unlike in a designed experiment, the probability of the instrument (prescribed therapy) conditional on the vector of prognostic factors would not be known and would have to be estimated from the data using a dimension reducing, usually parametric, model. In such cases, no estimators satisfying (i) and (ii) will exist. The best that can be hoped for with a logistic structural mean model is, as argued in this presentation, to succeed in constructing an estimator that is so-called doubly robust or doubly protected under the null hypothesis of no effect of either prescribed or actual treatment. For additive or multiplicative structural mean models, estimators that are also doubly robust at non-null parameter values are available (Robins, 2000).

We now summarize the model assumptions and the inferential problem treated in this presentation. We assume that we observe n iid copies O_i ; $i = 1, \dots, n$; of the vector $O = (L; Z; X; Y)$ where L is a vector recording pre-randomization covariates, Z records the randomly assigned treatment, X is the treatment actually received, and Y is the dichotomous indicator outcome of interest which we assume takes the values 0 or 1. Both Z and X can be discrete or continuous variables. Throughout, random variables are denoted with capital letters and their supports with the corresponding calligraphic letter.

To define treatment effects and our model we first define for each value $z \in Z$; $x \in X$ the potential outcome Y_{zx} (Neyman, 1923, Rubin, 1978, Robins 1986): This is defined as the value of the dichotomous outcome of interest had, possibly contrary to fact, Z been set to z and X set to x : Our semiparametric model, which throughout we denote by A ; has observed data O and is defined by the following restrictions: a) $f(z; L)$ is known by design, b) $Z = Y_{zx} \mid L$; w.p.1 for all $z \in Z$; $x \in X$; c) $Y = Y_{zx}$ if $Z = z$; $X = x$; d)

$$\theta(l; z; x) = \psi(\cdot; z; x; \tilde{A}_0) \tag{1}$$

where

$$\theta(z; x; l) = \log \frac{P(Y_{zx} = 1 \mid Z = z; X = x; L = l)^{1/2}}{P(Y_{zx} = 0 \mid Z = z; X = x; L = l)^{1/2}} \tilde{A} \frac{P(Y_{z0} = 1 \mid Z = z; X = x; L = l)^{3/4}}{P(Y_{z0} = 0 \mid Z = z; X = x; L = l)^{3/4}}$$

and \tilde{A}_0 is an unknown parameter vector of dimension p and $\psi(\cdot; z; x; \tilde{A})$ is a known function

satisfying $\phi(\cdot; z; 0; \bar{A}) = \phi(\cdot; z; x; 0) = 0$ and e) $b(z_0 j \cdot) = t(\cdot) + t^a(\cdot; z)$; with $t(\cdot)$ a completely unknown function, $t^a(\cdot; z)$ a known function satisfying $t(\cdot; 0) = 0$ and $b(zxj d) \sim \odot^{-1} fP(Y_{zx} = 1 | D = d)g$ where $\odot(u) = e^u = (1 + e^u)$ and $\odot^{-1}(u) = \log fu = (1 + u)g$: Our goal is to estimate the function $\phi(z; x; l)$ which measures the causal effect, on the logistic scale, of assigned treatment z and actual treatment x on the population assigned to z and treated with x and with baseline covariates equal to l : Randomization guarantees that Assumption b) holds. Assumption c) links the counterfactual data to the observed data. We make assumption d) in order to identify $\phi(z; x; l)$: We do so because without modeling restrictions on $\phi(l; z; x)$; this function is not generally identified from the observed data O : (See Robins, 1994). Following Robins (1997) and Robins et al (1999) we call the model (1); a logistic structural mean model. Note that $\phi(\cdot; z; x; 0) = 0$ implies that if $\bar{A}_0 = 0$ then the treatment $X = x$ has no effect on the outcome in the subpopulation with X observed to be x : Finally, assumption e) says that the direct effect of assigned treatment on the outcome when treatment actually taken is set to zero is known. Note that the special choice $t^a(\cdot; z) = 0$ is equivalent to the assumption that $P(Y_{z0} = 1 | L)$ is the same for all z and therefore is implied by (but does not imply) the exclusion assumption $Y_{zx} = Y_x$ w.p.1 for all $z \in Z$ and $x \in X$.

Though ideally we would like to make inference about \bar{A}_0 under model A; unfortunately, due to the curse of dimensionality (Huber, 1985, Robins and Ritov, 1997) this is not feasible unless $\bar{A}_0 = 0$ and $t^a(\cdot; z) = 0$. Thus, in order to obtain well behaved estimators of \bar{A}_0 when $\bar{A}_0 \neq 0$ we propose estimating \bar{A}_0 under a parametric model for the law of the observed data. However, we do not fit the model by maximum likelihood because the MLE of \bar{A}_0 is not guaranteed to either (i) consistently estimate the (null) treatment effect $\bar{A}_0 = 0$ when the sharp null hypothesis $Y_{zx} = Y_x$ w.p.1 of no effect of actual or assigned treatment is true or (ii) to have small bias when the true treatment effect \bar{A}_0 is close but not equal to zero and assigned treatment has no direct effect on the outcome i.e. $t^a(l; z) = 0$.

We now describe estimating equations whose solutions satisfy (i) and (ii). Let $q(\cdot; z; x) \sim b(z_0 j \cdot; z; x) - b(z_0 j \cdot; z; x = 0)$ where $b(zx^0 j l; z; x)$ denotes $\odot^{-1} fE(Y_{zx^0} | l; z; x)g$ and $\odot(u) := e^u = (1 + e^u)$: Robins and Rotnitzky (2001) showed that the observed data law can be parameterized in terms of $\phi(\cdot; z; x)$; $f(x j \cdot; z)$; $t(\cdot)$; and $q(\cdot; z; x)$ because $b(z_0 j \cdot; z; 0)$ depends on $q(\cdot; z; x)$; $t(\cdot)$ and $f(x j \cdot; z)$ since it is the unique solution to the integral equation

$$\odot f t(\cdot) + t^a(\cdot; z)g = \int \odot f q(\cdot; z; x) + v(\cdot; z)g dF(x j \cdot; z) \quad (2)$$

Robins and Rotnitzky (2001) provide the explicit analytical solution of (2) when X is binary.

Our proposed estimator of \bar{A} is computed by the following two-stage estimation procedure. At stage 1, postulate working parametric models $f(x j \cdot; z; \theta)$; $t(\cdot; !)$ and $q(\cdot; z; x; \mu)$ for the unknowns $f(x j \cdot; z)$; $t(\cdot)$ and $q(\cdot; z; x)$: For each fixed value of \bar{A} ; compute the estimator $\hat{b}(\bar{A}) = \hat{\theta}(\bar{A})$; $\hat{b}(\bar{A})$; $\hat{\mu}(\bar{A})$ of $\theta = (\theta; !; \mu)$; that maximize the parametric log-likelihood $\sum_{i=1}^n P(Y_{ij} | L_i; Z_i; X_i; \bar{A}; \frac{1}{2})^{Y_i} f(X_i j L_i; Z_i; \theta)$; or otherwise compute $\hat{\theta}(\bar{A}) = \hat{\theta}$ by maximizing $\sum_{i=1}^n f(X_i j L_i; Z_i; \theta)$ and then compute $\hat{\mu}(\bar{A})$; $\hat{b}(\bar{A})$ by maximizing the parametric conditional likelihood $\sum_{i=1}^n P(Y_{ij} | L_i; Z_i; X_i; \bar{A}; \mu; !)^{Y_i}$ over μ ; ! with $\hat{\theta}$ and \bar{A} held fixed. At stage 2, compute the estimator \hat{A} of \bar{A} that solves $\sum_{i=1}^n [b_{eff}(L_i; Z_i; \bar{A}) - E f b_{eff}(L_i; Z_i; \bar{A}) | L_i] H_i(\bar{A}; \hat{b}(\bar{A})) = 0$; where

$$H(\bar{A}; \frac{1}{2}) = \frac{\sum_{i=1}^n \odot [M_1(\frac{1}{2}) | L_i; Z_i] f Y_i - \odot [M_2(\bar{A}; \frac{1}{2})]g}{E_{\frac{1}{2}} f \odot [M_2(\bar{A}; \frac{1}{2})] | L_i; Z_i} + \sum_{i=1}^n f M_1(\frac{1}{2})g - \sum_{i=1}^n f M_3(\frac{1}{2})g$$

with $M_1(\frac{1}{2}) = q(L; Z; X; \mu) + v(L; Z; X; \frac{1}{2})$; $M_2(\tilde{A}; \frac{1}{2}) = \phi(L; Z; X; \tilde{A}) + q(L; Z; X; \mu) + v(L; Z; X; \frac{1}{2})$; $M_3(\frac{1}{2}) = t(L; !) + t^a(L; Z)$ and $c_{eff}(L_i; Z_i; \tilde{A}) = c_{eff}(L_i; Z_i; \tilde{A}; \frac{1}{2}(\tilde{A}))$ where

$$c_{eff}(L; Z; \tilde{A}; \frac{1}{2}) = \frac{E \left[s(L; Z; \tilde{A}; \frac{1}{2}) \sum_{i=1}^n \left[\frac{\partial}{\partial \tilde{A}} M_1(\frac{1}{2}) \right] \left[\frac{\partial}{\partial \tilde{A}} \phi(L; Z; X; \tilde{A}) \right] \right]}{E \left[s(L; Z; \tilde{A}; \frac{1}{2}) \sum_{i=1}^n \left[\frac{\partial}{\partial \tilde{A}} M_3(\frac{1}{2}) \right] \right]} s(L; Z; \tilde{A}; \frac{1}{2}) \sum_{i=1}^n \left[\frac{\partial}{\partial \tilde{A}} M_2(\frac{1}{2}) \right]$$

with $s(L; Z; \tilde{A}; \frac{1}{2}) = E \left[\frac{\partial}{\partial \tilde{A}} M_2(\frac{1}{2}) \right] \left[\frac{\partial}{\partial \tilde{A}} \phi(L; Z; X; \tilde{A}) \right]$: The asymptotic variance of $\hat{\tilde{A}}$ is equal to the semiparametric variance bound for \tilde{A} in model A under correct specification of the working models for $f(x; z)$; $t(\cdot)$ and $q(\cdot; z; x)$:

REFERENCES

- Abadie, A. (1999). Semiparametric instrumental variable estimation of treatment response models. Technical Report, MIT.
- Angrist, J. (2001). Estimation of Limited-dependent Variable Models with Dummy Endogenous Regressors: Simple Strategies for Empirical Practice. *J. of Business and Economic Statistics*, 19, 2-16.
- Hirano, K., Imbens, G., Rubin, D., and Zhou, A. (2000). Assessing the effect of an influenza vaccine in an encouragement design. *Biostatistics*, 1.1, 69-88.
- Huber, P. (1985). Projection Pursuit. *Annals of Statistics*, 13, 435-475.
- Neyman, J. (1923). On the application of probability theory to agricultural experiments: essay on principles, section 9. Translated in *Statistical Science*, 5, 565-480, 1990
- Robins, J.M. (1986). 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.
- Robins, J. M. (1994) Correcting for non-compliance in randomized trials using structural nested mean models, *Communications in Statistics, Theory and Methods*, 23, 2379-2412.
- Robins, J.M. (1997). Causal Inference from Complex Longitudinal Data. *Latent Variable Modeling and Applications to Causality*. Lecture Notes in Statistics (120), M. Berkane, Editor. NY: Springer Verlag, pp. 69-117.
- Robins, J. and Ritov, Y. (1997). Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semiparametric models. *Statistics in Medicine*. 16, 285-319.
- Robins, J.M. (1999). Marginal Structural models vs structural nested mean models as tools for causal inference. *Statistical Models in Epidemiology: The environment and Clinical Trials*, New York: Springer Verlag, 95-134.
- Robins, J.M. (2000). Robust estimation in sequentially ignorable missing data and causal inference models. *Proceedings of the American Statistical Association Section on Bayesian Statistical Science 1999*, pp. 6-10.
- Robins, J.M. and Rotnitzky, A. (2001). Estimation of treatment effects in randomized trials with noncompliance and a dichotomous outcome using logistic structural mean models. Submitted.
- Rubin, D. (1978). Bayesian inference for causal models. *Annals of Statistics*, 6, 34-58.
- Vansteelandt, S. and Goetghebeur, E. (2001) Causal inference with generalized structural mean models. To appear in the *Journal of the Royal Statistical Society, series B*.