

a summary of known results concerning the existence and construction of doubly robust estimators, including results for missing data models with data that are not CAR.

One could wonder about the actual advantage of using doubly robust estimators as, in practice, all models including the lower dimensional models for both  $G$  and  $F_X$  are misspecified. Thus, even a doubly robust estimator of  $\mu$  may be considerably biased. In our opinion, however, a doubly robust estimator has the following advantage that argues for its routine use: if either the lower dimensional model for  $F_X$  or for  $G$  is nearly correct, then the bias of doubly robust estimator of  $\mu$  will be small. Thus, a doubly robust estimator gives the analyst two chances, instead of only one, to obtain a nearly unbiased estimator of  $\mu$ . Furthermore informal but sensitive goodness of fit tests can be based on doubly robust estimators. See Robins and Rotnitzky (2001) for details. Yu, van der Laan (2002) demonstrate in a simulation study that the maximum likelihood estimator is much less robust than doubly robust estimators to model misspecification.

Van der Vaart, van der Laan (2001), van der Laan, Yu (2001), and Robins and Rotnitzky (2001) show that our general estimating function methodology can also be used to obtain doubly robust estimators of non-smooth parameters (such as the density of a continuous random variable) in both CAR censored data models and in certain other semiparametric models by approximating the non-smooth parameter by a smooth parameter, applying the locally efficient estimating function methodology to estimate the smooth parameter, and allowing the approximation to improve at an appropriate rate with increasing sample size. In van der Vaart, van der Laan (2001) this approach is used to estimate a survival function based on current status data in the presence of a time-dependent surrogate processes. In this monograph, however, we largely restrict attention to the estimation of smooth parameters.

In this book we assume that the data are CAR. This assumption will in general be correct only when data have been obtained on all time-independent and time-dependent covariates that predict both (i) subsequent outcomes and (ii) subsequent treatment and/or censoring. In certain settings CAR may not hold even approximately. For example subjects often refrain from answering specific questions in surveys of religious, financial or sexual practices, but their reasons for nonresponse (censoring) are unavailable to the data analyst. Thus methods for the analysis of non-CAR coarsened data are also important. In a series of papers, Robins, Rotnitzky and Scharfstein have developed a general estimating function sensitivity analysis methodology for locally semiparametric efficient (LSE) estimation of smooth parameters (i.e., parameters estimable at rate  $\sqrt{n}$ )  $\mu$  of large non- or semi- parametric models for the law  $F_X$  of very high dimensional data  $X$  based on the observed data  $Y$ , when the data are not coarsened at random. For further information, the interested reader may consult Rotnitzky, Robins (1997), Rotnitzky, Robins, and Scharfstein

(1998) Scharfstein, Rotnitzky, and Robins (1999), and Robins, Rotnitzky, Scharfstein (1999).

Robins (1987), Robins and Rotnitzky (1992), Robins, Rotnitzky, and Zhao (1995), Robins (1997), Pearl (1995), and Pearl and Robins (1995) provide conditions under which CAR is violated but causal effects can still be nonparametrically identified by the  $G$ -computation functional algorithm of Robins (1986). Furthermore the induced model for the observable data  $Y$  under these conditions is the same as under CAR, so the theory and methodology developed in this book for CAR models still applies. See Subsection 3.3.2 and Section 6.2 for examples. Pearl (1995,2000) describes additional assumptions characterized in terms of missing arrows on "causal" directed acyclic graphs under which causal effects are nonparametrically identified by a functionals of the observed data distribution other than the  $G$ -computation algorithm functional. Though these additional identification results are quite interesting, the application of these latter results to complex longitudinal studies with high dimensional time-dependent data structures as considered in this book has been limited because their use requires detailed knowledge of the underlying causal structure, which rarely exists.

The book is organized as follows. Chapter 1 provides an overview of the mathematical tools that lie at the foundation of our methodology. Specifically we review the modern theory of semiparametric models based on tangent spaces and Hilbert space projections on these spaces. We then use this theory to derive our estimating function methodology and prove that these estimating functions are doubly robust. Next we apply our general estimating function methodology to several important examples. Chapter 2 is more technical and provides a rather complete treatment of our locally efficient estimating function methodology and its associated asymptotic theory. Although the theory and methods developed in this chapter are illustrated with concrete examples, nonetheless Chapter 2 is the most difficult in the book. The remaining chapters 3-6, however, cover specific data structures and models and are sufficiently self-contained that they can be read before chapter 2, although occasionally one will be referred back to a result in chapter 2. Chapter 3 treats right-censored data. Chapter 4 treats current status data and a combination of current status data and right-censored data. Chapter 5 treats multivariate right-censored data. Finally, Chapter 6 provides a general methodology for estimation of causal and non-causal parameters from longitudinal data, complicated by right censoring of some responses and interval censoring of other responses. Most of the methods discussed in chapters 3-6 have been implemented with Splus or C-software by Ph.D students in Biostatistics and Statistics. Examples of these data analyses and simulations are provided toward the end of each chapter.

With the help of students and colleagues, we have applied our locally efficient, doubly robust, estimating function methodology to a large variety