

NBER WORKING PAPER SERIES

EFFICIENT ESTIMATION OF MISSING DATA MODELS USING MOMENT CONDITIONS
AND SEMIPARAMETRIC RESTRICTIONS

Bryan S. Graham

Working Paper 14376
<http://www.nber.org/papers/w14376>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
October 2008

I would like to thank Gary Chamberlain, Jinyong Hahn, Guido Imbens, Michael Jansson and Whitney Newey for comments on earlier draft. Helpful discussions with Oliver Linton, Cristine Pintos, Jim Powell, Geert Ridder as well as participants in the Berkeley Econometrics Reading Group and Seminars are gratefully acknowledged. This revision has benefited from Tom Rothenberg's skepticism, discussions with Michael Jansson, Justin McCrary, Jim Powell, the comments of a co-editor and two anonymous referees. All the usual disclaimers apply. This is a heavily revised version of a paper which previously circulated under the titles "A note on semiparametric efficiency in moment condition models with missing data" and "GMM 'equivalence' for semiparametric missing data models.". The views expressed herein are those of the author(s) and do not necessarily reflect the views of the National Bureau of Economic Research.

© 2008 by Bryan S. Graham. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Efficient Estimation of Missing Data Models Using Moment Conditions and Semiparametric Restrictions

Bryan S. Graham

NBER Working Paper No. 14376

October 2008

JEL No. C1,C14,C21

ABSTRACT

This paper shows that the semiparametric efficiency bound for a parameter identified by an unconditional moment restriction with data missing at random (MAR) coincides with that of a particular augmented moment condition problem. The augmented system consists of the inverse probability weighted (IPW) original moment restriction and an additional conditional moment restriction which exhausts all other implications of the MAR assumption. The paper also investigates the value of additional semiparametric restrictions on the conditional expectation function (CEF) of the original moment function given always-observed covariates. In the missing outcome context, for example, such restrictions are implied by a semiparametric model for the outcome CEF given always-observed covariates. The efficiency bound associated with this model is shown to also coincide with that of a particular moment condition problem. Some implications of these results for estimation are briefly discussed.

Bryan S. Graham

UC, Berkeley

Department of Economics

508-1 Evans Hall #3880

Berkeley, CA 94720-3880

and NBER

bgraham@econ.berkeley.edu

1 Introduction

Let $Z = (Y_1', X)'$ be vector of modelling variables, $\{Z_i\}_{i=1}^\infty$ be an independent and identically distributed random sequence drawn from the unknown distribution F_0 , β a $K \times 1$ unknown parameter vector and $\psi(Z, \beta)$ a known vector-valued function of the same dimension.² The only prior restriction on F_0 is that for some $\beta_0 \in \mathcal{B} \subset \mathbb{R}^K$

$$\mathbb{E}[\psi(Z, \beta_0)] = 0. \quad (1)$$

Chamberlain (1987) showed that the maximal asymptotic precision with which β_0 can be estimated under (1) (subject to identification and regularity conditions) is given by the inverse of $\mathcal{I}_f(\beta_0) = \Gamma_0' \Omega_0^{-1} \Gamma_0$, with $\Gamma_0 = \mathbb{E}[\partial \psi(Z, \beta_0) / \partial \beta']$ and $\Omega_0 = \mathbb{V}(\psi(Z, \beta_0))$.³

Now consider the case where a random sequence from F_0 is unavailable. Instead only a ‘selected’ sequence of samples is available. Let D be a binary selection indicator. When $D = 1$ we observe Y_1 and X , when $D = 0$ we observe only X .⁴ This paper considers estimation of β_0 under restriction (1) and the following additional assumptions.

Assumption 1.1 (RANDOM SAMPLING) $\{Z_i, D_i\}_{i=1}^\infty$ is an independent and identically distributed random sequence from F_0 .

Assumption 1.2 (OBSERVED DATA) For each unit we observe X , D and $Y = DY_1$.

Assumption 1.3 (CONDITIONAL INDEPENDENCE) $Y_1 \perp D | X$.

Assumption 1.4 (OVERLAP) Let $p_0(x) = \Pr(D = 1 | X = x)$, then $0 < \kappa \leq p_0(x) \leq 1$ for all $x \in \mathcal{X} \subset \mathbb{R}^{\dim(x)}$.

Restriction (1) and Assumptions 1.1 to 1.4 constitute a semiparametric model for the data. Henceforth I refer to this model as the semiparametric missing data model or the missing at random (MAR) setup. Robins, Rotnitzky and Zhao (1994, Proposition 2.3, p. 850) derived the efficient influence function for this problem and proposed a locally efficient augmented inverse probability weighting (AIPW) estimator (cf., Scharfstein, Rotnitzky and Robins 1999, Bang and Robins 2005,

²Extending what follows to the overidentified case is straightforward.

³Throughout upper case letters denote random variables, lower case letters specific realizations of them, and calligraphic letters their support. I use the notation $\mathbb{E}[A|c] = \mathbb{E}[A|C=c]$, $\mathbb{V}(A|c) = \text{Var}(A|C=c)$ and $\mathbb{C}(A, B|c) = \text{Cov}(A, B|C=c)$.

⁴An earlier version of this paper considered the slightly more general set-up with $\psi(Z, \beta) = \psi_1(Y_1, X, \beta) - \psi_0(Y_0, X, \beta)$ with (X, Y) observed where $Y = DY_1 + (1 - D)Y_0$. Results for this extended model, which contains the standard causal inference model and the two-sample instrumental variables model as special cases (cf., Imbens 2004, Angrist and Krueger 1992), follow directly and straightforwardly from those outlined below.

Tsiatis 2006). Cheng (1994), Hahn (1998), Hirano, Imbens and Ridder (2003), Imbens, Newey and Ridder (2005), and Chen, Hong and Tarozzi (2008) develop globally efficient estimators.

The ‘MAR setup’ has been applied to a number of important econometric and statistical problems, including program evaluation as surveyed by Imbens (2004), non-classical measurement error (e.g., Robins, Hsieh and Newey 1995, Chen, Hong and Tamer 2005), missing regressors (e.g., Robins, Rotnitzky and Zhao 1994), attrition in panel data (e.g., Wooldridge 2002), and M-estimation under variable probability sampling (e.g., Wooldridge 1999, 2007). Chen, Hong and Tarozzi (2004), Wooldridge (2007) and Egel, Graham and Pinto (2008) discuss several other applications.

The maximal asymptotic precision with which β_0 can be estimated under the MAR setup has been characterized by Robins, Rotnitzky and Zhao (1994) and is given by the inverse of

$$\mathcal{I}_m(\beta_0) = \Gamma'_0 \mathbb{E}[\Lambda_0(X)]^{-1} \Gamma_0, \quad (2)$$

with $\Lambda_0(x) = \Sigma_0(x)/p_0(x) + q(x; \beta_0)q(x; \beta_0)'$, where $\Sigma_0(x) = \mathbb{V}(\psi(Z, \beta_0)|x)$ and $q(x; \beta) = \mathbb{E}[\psi(Z, \beta)|x]$.

The associated efficient influence function, also due to Robins, Rotnitzky and Zhao (1994), is given by

$$\phi(z, \theta_0) = \Gamma_0^{-1} \times \left\{ \frac{d}{p_0(x)} \psi_1(z, \beta_0) - \frac{q(x; \beta_0)}{p_0(x)} (d - p_0(x)) \right\} \quad (3)$$

for $\theta = (p, q', \beta')'$.

The calculation of (2) is now standard. Knowledge of (2) is useful because it quantifies the cost – in terms of asymptotic precision – of the missing data and because it can be used to verify whether a specific estimator for β_0 is efficient. To simplify what follows I will explicitly assume that $\mathcal{I}_m(\beta_0)$ is well-defined (i.e., that all its component expectations exist and are finite, and that all its component matrices are nonsingular).

This paper shows that the semiparametric efficiency bound for β_0 under the MAR setup, coincides with the bound for a particular augmented moment condition problem. The augmented system consists of the inverse probability of observation weighted (IPW) original moment restriction (1) and an additional conditional moment restriction which exhausts all other implications of the MAR setup (useful for estimating β_0). This general equivalence result, while implicit in the form of the efficient influence function (3), is apparently new. It provides fresh intuitions for several ‘paradoxes’ in the missing data literature, including the well-known results that projection onto, or weighting by the inverse of, a known propensity score results in inefficient estimates (e.g., Hahn 1998, Hirano, Imbens and Ridder 2003), that smoothness and exclusion priors on the propensity score do not increase the precision with which β_0 can be estimated (Robins, Hsieh and Newey 1995, Robins and Rotnitzky 1995, Hahn 1998, 2004) and that weighting by a nonparametric estimate of the propensity

score results in an efficient estimator (Hirano, Imbens and Ridder 2003, cf., Hahn 1998, Wooldridge 2007).

This paper also analyzes the effect of imposing additional semiparametric restrictions on the conditional expectation function (CEF) $q(x; \beta) = \mathbb{E}[\psi(Z, \beta) | x]$. If $\psi(Z, \beta) = Y_1 - \beta$, as when the target parameter is $\beta_0 = \mathbb{E}[Y_1]$, then such restrictions may arise from prior information on the form of $\mathbb{E}[Y_1 | x]$. Such restrictions may arise in other settings as well. For example, if the goal is to estimate a vector of linear predictor coefficients in the presence of missing regressors, then a semiparametric model for the CEFs of the missing regressors given always-observed variables generates restrictions on the form of $q(x; \beta)$ (cf., Robins, Rotnitzky and Zhao 1994).⁵

Formally I consider the semiparametric model defined by restriction (1), Assumptions 1.1 to 1.4 and the additional assumption.

Assumption 1.5 (FUNCTIONAL RESTRICTION) *Partition $X = (X'_1, X'_2)'$, then*

$$\mathbb{E}[\psi(Z, \beta_0) | x] = q(x, \delta_0, h_0(x_2); \beta_0)$$

where $q(x, \delta, h(x_2); \beta)$ is a known $K \times 1$ function, δ a $J \times 1$ finite dimensional unknown parameter, and $h(\cdot)$ an unknown function mapping from a subset of $\mathcal{X}_2 \subset \mathbb{R}^{\dim(X_2)}$ into $\mathcal{H} \subset \mathbb{R}^P$.

To the best of my knowledge the variance bound for this problem, the MAR setup with ‘functional’ restrictions, has not been previously calculated. In an innovative paper, Wang, Linton and Härdle (2004) consider a special case of this model where $\psi(Z, \beta) = Y_1 - \beta$. They impose a partial linear structure, as in Engle et al (1986), on $\mathbb{E}[Y_1 | x]$ such that $q(x, \delta_0, h_0(x_2); \beta_0) = x'_1 \delta_0 + h_0(x_2) - \beta_0$. In making their variance bound calculation they assume that the conditional distribution of Y_1 given X is normal with a variance that does not depend on X . They do not provide a bound for the general case but conjecture that it is “very complicated” (p. 338). The result given below extends their work to moment condition models, general forms for $q(x, \delta, h(x_2); \beta)$ and, importantly, does not require that $\psi(Z, \beta)$ be conditionally normally distributed and/or homoscedastic.

Augmenting the MAR setup with Assumption 1.5 generates a middle ground between the fully parametric likelihood-based approaches to missing data described by Little and Rubin (2002) and those which leave $\mathbb{E}[\psi(Z, \beta_0) | x]$ unrestricted (e.g., Cheng 1994, Hahn 1998, Hirano, Imbens and Ridder 2003). Likelihood-based approaches are very sensitive to misspecification (cf., Imbens 2004), while approaches which utilize only the basic MAR setup require high dimensional smoothing which may deleteriously affect small sample performance (cf., Wang, Linton and Härdle 2004, Ichimura and Linton 2005). Assumption 1.5 is generally weaker than a parametric specification for the

⁵The formation of predictive models of this type is the foundation of the imputation approach to missing data described in Little and Rubin (2002).

conditional distribution of $\psi(Z, \beta_0)$ given X , but at the same time reduces the dimension of the nonparametric smoothing problem. Below I show how to efficiently exploit prior information on the form of $\mathbb{E}[\psi(Z, \beta_0) | x]$. I also provide conditions under which consistent estimation of β_0 is possible even if the exploited information is incorrect.

Section 2 reports the first result of the paper: an equivalence between the ‘MAR setup’ and a particular method-of-moments problem. Equivalence, which is suggested by the form of the efficient influence function derived by Robins, Rotnitzky and Zhao (1994), was previously noted for special cases by Newey (1994a) and Hirano, Imbens and Ridder (2003). I discuss the connection between their results and the general result provided below. I also highlight some implications of the equivalence result for understanding various aspects of the MAR setup. Section 3 calculates the variance bound for β_0 when the MAR setup is augmented by Assumption 1.5. I discuss when Assumption 1.5 is likely to be informative and also when consistent estimation is possible even if it is erroneously maintained.

2 Equivalence result

Under the MAR setup the inverse probability weighted (IPW) moment condition

$$\mathbb{E} \left[\frac{D}{p_0(X)} \psi(Z, \beta_0) \right] = 0, \quad (4)$$

is valid (e.g., Hirano, Imbens and Ridder 2003, Wooldridge 2007). The conditional moment restriction

$$\mathbb{E} \left[\frac{D}{p_0(X)} - 1 \middle| X \right] = 0 \quad \forall \quad X \in \mathcal{X}, \quad (5)$$

also holds and nonparametrically identifies $p_0(x)$. While the terminology is inexact, in what follows I call (4) the *identifying moment* and (5) the *auxiliary moment*.

Consider the case where $p_0(x)$ is known such that (5) is truly an auxiliary moment. One efficient way to exploit the information (5) contains is to, following Newey (1994a) and Brown and Newey (1998), reduce the sampling variation in (4) by subtracting from it the fitted value associated with its regression onto the infinite-dimensional vector of unconditional moment functions implied by (5):

$$\begin{aligned} s(Z, \theta_0) &= \frac{D}{p_0(X)} \psi(Z, \beta_0) - \mathbb{E}^* \left[\frac{D}{p_0(X)} \psi(Z, \beta_0) \middle| \frac{D}{p_0(X)} - 1; X \right] \\ &= \frac{D}{p_0(X)} \psi(Z, \beta_0) - \frac{q(X; \beta_0)}{p_0(X)} (D - p_0(X)). \end{aligned}$$

That this ‘population residual’ is equal to the efficient score function derived by Robins, Rotnitzky

and Zhao (1994) strongly suggests an equivalence between the GMM problem defined by restrictions (4) and (5) and the MAR setup outlined above.⁶ That this is indeed the case is shown by the following Theorem.⁷

Theorem 2.1 (GMM EQUIVALENCE) *Suppose that (i) the distribution of Z has a known, finite support, (ii) there is some $\beta_0 \in \mathcal{B} \subset \mathbb{R}^K$ and $\rho_0 = (\rho_1, \dots, \rho_L)'$ where $\rho_l = p_0(x_l) \in \mathcal{P} \subset [\kappa, 1]$ for each $l = 1, \dots, L$ and some $0 < \kappa < 1$ (with $\mathcal{X} = \{x_1, \dots, x_L\}$ the known support of X) such that restrictions (4) and (5) hold, (iii) $\mathbb{E}[\Lambda_0(X)]$ and $\mathcal{I}_m(\beta_0) = \Gamma_0' \mathbb{E}[\Lambda_0(X)]^{-1} \Gamma_0$ are nonsingular and (iv) other regularity conditions hold (cf., Chamberlain 1992b, Section 2), then $\mathcal{I}_m(\beta_0)$ is the Fisher information bound for β_0 .*

Proof. See Appendix A. ■

The proof of Theorem 2.1 involves only some tedious algebra and a straightforward application of Lemma 2 of Chamberlain (1987). Assuming that Z has known, finite support makes the problem fully parametric. The unknown parameters are the probabilities associated with each possible realization of Z , the values of the propensity score at each of the L mass points of the distribution of X , $\rho_0 = (\rho_1, \dots, \rho_L)'$, and the parameter of interest, β_0 .

The multinomial assumption is not apparent in the form of $\mathcal{I}_m(\beta_0)$, which involves only conditional expectations of certain functions of the data. This suggests that the bound holds in general since any F_0 which satisfies (4) and (5) can be arbitrarily well-approximated by a multinomial distribution also satisfying the restrictions. Chamberlain (1992a, Theorem 1) demonstrates that this is indeed the case. Therefore $\mathcal{I}_m(\beta_0)^{-1}$ is the maximal asymptotic precision, in the sense of Hájek's (1972) local minimax approach to efficiency, with which β_0 can be estimated when the only prior restrictions on F_0 are (4) and (5). Since this variance bound coincides with (2) I conclude that (4) and (5) 'exhaust' all of the prior restrictions implied by the MAR setup (that are helpful for estimating β_0).

The connection between semiparametrically efficient estimation of moment condition models with missing data and 'augmented' systems of moment restrictions has been noted previously for the special case of data missing completely at random (MCAR). In that case Assumptions 1.1 to 1.4 hold with $p_0(X)$ equal to a (perhaps known) constant. Newey (1994a) shows that an efficient estimate of β_0 can be based on the pair of moment restrictions

$$\mathbb{E}[D\psi(Z, \beta_0)] = 0, \quad \mathbb{C}(D, q(X; \beta_0)) = 0,$$

⁶The notation $\mathbb{E}^*[Y|X; Z]$ denotes the (mean squared error minimizing) linear predictor of Y given X within a subpopulation homogenous in Z . Wooldridge (1999b, Section 4) collects some useful results on conditional linear predictors. See also Newey (1990) and Brown and Newey (1998).

⁷An alternative approach to showing equivalency would involve verifying Newey's (2004) moment spanning condition for efficiency.

with $q(X; \beta)$ as defined above. Hirano, Imbens and Ridder (2003) discuss a related example with X binary and the data also MCAR. In their example efficient estimation is possible with only a finite number of unconditional moment restrictions. Theorem 2.1 provides a formal generalization of the Newey (1994a) and Hirano, Imbens and Ridder (2003) examples to the missing at random (MAR) case.

The method-of-moments formulation of the MAR setup provides a useful framework for understanding several apparent paradoxes found in the missing data literature. As a simple example consider Hahn's (1998, pp. 324 - 325) result that projection onto a known propensity score may be harmful for estimation of $\beta_0 = \mathbb{E}[Y_1]$. Formally he shows that, for $p_0(x) = Q_0$ constant in x and known, the complete-case estimator

$$\hat{\beta}_{cc} = \sum_{i=1}^N D_i Y_{1i} / \sum_{i=1}^N D_i,$$

while consistent, is inefficient. Observe that for the constant propensity score case $\hat{\beta}_{cc}$ is the sample analog of the population solution to (4). It consequently makes no use of any information contained in the auxiliary moment (5). However, that moment will be informative for β_0 if $q(x; \beta_0) = \mathbb{E}[Y_1 | x] - \beta_0$ varies with x , consistent with Hahn's (1998) finding that the efficiency loss associated with $\hat{\beta}_{cc}$ is proportional to $\mathbb{V}(q(X; \beta_0))$. Similar reasoning explains why weighting by the (inverse of) the known propensity score is generally inefficient (cf., Robins, Rotnitzky and Zhao 1994, Hirano, Imbens and Ridder 2003, Wooldridge 2007). The known weights estimator ignores the information contained in (5).

That smoothness and exclusion priors on the propensity score do not lower the variance bound also has a GMM interpretation. Consider the case where the propensity score belongs to a parametric family $p(X; \eta_0)$. If η_0 is known, then an efficient GMM estimator based on (4) and (5) is given by the solution to

$$\frac{1}{N} \sum_{i=1}^N s(\eta_0, \hat{q}, \hat{\beta}) = \frac{1}{N} \sum_{i=1}^N \left\{ \frac{D_i}{p(X_i; \eta_0)} \psi(Z_i, \hat{\beta}) - \frac{\hat{q}(X_i; \hat{\beta})}{p(X_i; \eta_0)} (D_i - p(X_i; \eta_0)) \right\} = 0,$$

with $\hat{q}(x; \hat{\beta})$ a consistent nonparametric estimate of $\mathbb{E}[\psi(Z, \beta_0) | x]$. Now consider the effect of replacing η_0 with the consistent estimate $\hat{\eta}$. From Newey and McFadden (1994, Theorem 6.2), this replacement does not change the first order asymptotic sampling distribution of $\hat{\beta}$ because $\mathbb{E}[\partial s(\eta_0, q_0, \beta_0) / \partial \eta'] = 0$. Furthermore, if the known propensity score is replaced by a consistent nonparametric estimate, $\hat{p}(x)$, then the sampling distribution of $\hat{\beta}$ is also unaffected (Newey 1994b, Proposition 3, p. 1360). Since the M-estimate of β_0 based on its efficient score function has the same

asymptotic sampling distribution whether the propensity score is set equal to the truth or instead to a noisy, but consistent, estimate, knowledge of its form cannot increase the precision with which β_0 may be estimated.

Another intuition for redundancy of knowledge of the propensity score can be found by inspecting the information bound for the multinomial problem. Under the conditions of Theorem 2.1 calculations provided in Appendix A imply that the GMM estimates of β_0 and ρ_0 (recall that ρ_0 contains the values for the propensity score at each of the mass points of the distribution of X) have an asymptotic sampling distribution of

$$\sqrt{N} \left(\begin{bmatrix} \hat{\rho} \\ \hat{\beta} \end{bmatrix} - \begin{bmatrix} \rho_0 \\ \beta_0 \end{bmatrix} \right) \xrightarrow{D} \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \mathcal{I}_m(\rho_0)^{-1} & 0 \\ 0 & \mathcal{I}_m(\beta_0)^{-1} \end{bmatrix} \right),$$

with $\mathcal{I}_m(\beta_0)$ as defined in (2) and $\mathcal{I}_m(\rho_0)$ as defined in Appendix A. As is well-known, under block diagonality sampling error in $\hat{\rho}$ does not affect, at least to first order, the asymptotic sampling properties of $\hat{\beta}$. While block diagonality is formally only a feature of the multinomial problem, the result nonetheless provides another useful intuition for understanding why prior knowledge of the propensity score is not valuable asymptotically.

Finally the combination of redundancy of knowledge of the propensity score, and the structure of the equivalent GMM problem, suggests why the IPW estimator based on a nonparametric estimate of the propensity score is semiparametrically efficient (Hirano, Imbens and Ridder 2003): when a nonparametric estimate of the propensity score is used the sample analog of both (4) and (5) are satisfied. In contrast the IPW estimator based on a parametric estimate of the propensity score will only satisfy a finite number of the moment conditions implied by (5), hence while it will be more efficient than the estimator which weights by the true propensity score (e.g., Wooldridge 2007), it will be less efficient than the one proposed by Hirano, Imbens and Ridder (2003).

3 Semiparametric functional restrictions

Consider the MAR setup augmented by Assumption 1.5. To the best of my knowledge, the maximal asymptotic precision with which β_0 can be estimated in this model has not been previously characterized. In order to calculate the bound for this problem I first consider the conditional moment problem defined by (4) and (5) and

$$\mathbb{E}[\rho(Z, \delta_0, h_0(X_2); \beta_0) | X] = 0. \quad (6)$$

with $\rho(Z, \delta_0, h_0(X_2); \beta_0) = \psi(Z, \beta_0) - q(x, \delta_0, h_0(x_2); \beta_0)$.

I apply Chamberlain's (1992a) approach to this problem to calculate a variance bound for β_0 . I

then verify that this bound is indeed the semiparametric efficiency bound for the problem defined by restriction (1) and Assumptions 1.1 to 1.5 using the methods of Bickel, Klaassen, Ritov and Wellner (1993). The value of first considering the conditional moment problem is that it provides a conjecture for the form of the efficient influence function, therefore sidestepping the need to directly calculate what is evidently a complicated projection. That Chamberlain's (1992a) bound is identical to the Bickel, Klaassen, Ritov and Wellner (1993) one also results in the paper's second GMM equivalence result.

To present these results I begin by letting $q_0(X) = q(X, \delta_0, h_0(X_2); \beta_0)$, $\rho(Z; \beta_0) = \psi(Z, \beta_0) - q_0(X)$, $\Sigma_0(X) = \mathbb{V}(\rho(Z; \beta_0) | X)$ and

$$\begin{aligned}\Upsilon_0^\delta(X_2) &= \mathbb{E} \left[D \left(\frac{\partial q_0(X)}{\partial \delta'} \right)' \Sigma_0(X)^{-1} \left(\frac{\partial q_0(X)}{\partial \delta'} \right) \middle| X_2 \right] \\ \Upsilon_0^h(X_2) &= \mathbb{E} \left[D \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right) \middle| X_2 \right] \\ \Upsilon_0^{h\delta}(X_2) &= \mathbb{E} \left[D \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \left(\frac{\partial q_0(X)}{\partial \delta'} \right) \middle| X_2 \right] \\ G_0(X) &= \left(\frac{\partial q_0(X)}{\partial \delta'} - \left(\frac{\partial q_0(X)}{\partial h'} \right) \Upsilon_0^h(X_2)^{-1} \Upsilon_0^{h\delta}(X_2) \right), \quad H_0(X_2) = \mathbb{E} \left[\frac{\partial q_0(X)}{\partial h'} \middle| X_2 \right] \\ \mathcal{I}_m^f(\delta_0) &= \mathbb{E} \left[D G_0(X)' \Sigma_0(X)^{-1} G_0(X) \right],\end{aligned}$$

and

$$\Xi_0 = \mathbb{E} \left[H_0(X_2) \Upsilon_0^h(X_2)^{-1} H_0(X_2)' \right] + \mathbb{E} [G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} \mathbb{E} [G_0(X)]' + \mathbb{E} [q_0(X) q_0(X)'] .$$

The variance bound for β_0 in the conditional moment problem defined by (4), (5) and (6) is established by the following Theorem.

Theorem 3.1 (EFFICIENCY WITH SEMIPARAMETRIC FUNCTIONAL RESTRICTIONS, PART 1) *Suppose that (i) the distribution of Z has a known, finite support, (ii) there is some $\beta_0 \in \mathcal{B} \subset \mathbb{R}^K$, $\rho_0 = (\rho_1, \dots, \rho_L)'$ where $\rho_l = p_0(x_l) \in \mathcal{P} \subset [\kappa, 1]$ for each $l = 1, \dots, L$ and some $0 < \kappa < 1$ (with $\mathcal{X} = \{x_1, \dots, x_L\}$ the known support of X), $\delta_0 \in \mathcal{D} \subset \mathbb{R}^J$ and $h_0(x_{2,m}) = \lambda_{0,m} \in \mathcal{L} \subset \mathbb{R}^P$ for each $m = 1, \dots, M$ (with $\mathcal{X}_2 = \{x_{2,1}, \dots, x_{2,M}\}$ the known support of X_2) such that restrictions (4), (5) and (6) hold, (iii) Ξ_0 and $\mathcal{I}_m^f(\beta_0) = \Gamma_0' \Xi_0^{-1} \Gamma_0$ are nonsingular and (iv) other regularity conditions hold (cf., Chamberlain 1992b, Section 2), then $\mathcal{I}_m^f(\beta_0)$ is the Fisher information bound for β_0 .*

Proof. See Appendix A. ■

The form of Ξ_0 suggests a candidate efficient influence function of

$$\begin{aligned} \phi_\beta^f(Z, \eta_0, \beta_0) = & \Gamma_0^{-1} \left\{ DH_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) \right. \\ & \left. + D\mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} G_0(X)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) + q(X; \beta_0) \right\}. \end{aligned} \quad (7)$$

with $\eta = (h, \delta, H, \Upsilon^h, G)$. Note that each of the three components of (7) are mutually uncorrelated. The next Theorem verifies this conjecture.

Theorem 3.2 (EFFICIENCY WITH SEMIPARAMETRIC FUNCTIONAL RESTRICTIONS, PART 2) *The semiparametric efficiency bound for β_0 in the problem defined by restriction (1) and Assumptions 1.1 to 1.5 is equal to $\mathcal{I}_m^f(\beta_0)$ with an efficient influence function of $\phi_\beta^f(Z, \eta_0, \beta_0)$.*

The forms of $\mathcal{I}_m^f(\beta_0)$ and $\phi_\beta^f(Z, \eta_0, \beta_0)$ provide insight into the costs and benefits of utilizing Assumption 1.5. If $X_2 = X$ such that $h_0(X_2) = q(X; \beta_0)$, then the first term in Ξ_0 is identically equal to $\Sigma_0(X)/p_0(X)$. Therefore, very loosely speaking, the second term in Ξ_0 represents the information gain associated with Assumption 1.5. From Chamberlain (1992a) the information bound for δ_0 in the semiparametric regression problem

$$D\psi(Z, \beta_0) = Dq(X, \delta_0, h_0(X_2); \beta_0) + DV, \quad \mathbb{E}[V|X, D=1] = \mathbb{E}[V|X] = 0,$$

is given by $\mathcal{I}_m^f(\delta_0)$. Therefore the more precisely identified δ_0 , the greater the efficiency gain from imposing Assumption 1.5. The form of $\mathbb{E}[G_0(X)]$ also governs the magnitude of any efficiency gain. Note that $\left(\frac{\partial q_0(X)}{\partial h'} \right) \Upsilon_0^h(X_2)^{-1} \Upsilon_0^{h\delta}(X_2) = \mathbb{E}_\Sigma^* \left[\frac{\partial q_0(X)}{\partial \delta'} \middle| \frac{\partial q_0(X)}{\partial h'}; X_2, D=1 \right]$,⁸ therefore $G_0(X)$ is equal to the difference between $\frac{\partial q_0(X)}{\partial \delta'}$ and its predicted value based on a weighted least squares regression in the $D=1$ subpopulation. The average of these differences, $\mathbb{E}[G_0(X)]$, will be large in absolute value when the distribution of X_1 conditional on X_2 differs in the $D=1$ versus $D=0$ subpopulations. This, in turn, will be true whenever X_1 is highly predictive for missingness conditional on X_2 . Intuitively in that case requiring that X_1 enter $q(X, \delta_0, h_0(X_2); \beta_0)$ parametrically mitigates the efficiency costs of (conditional) covariate imbalance across the two subpopulations by facilitating extrapolation.

The form of $\mathcal{I}_m^f(\beta_0)$ suggests when imposing Assumption 1.5 is likely to be informative. A related question concerns the consequences of misspecifying the form of $q(X, \delta, h(X_2); \beta)$. Under

⁸ $\mathbb{E}_\Sigma^*[Y|X; Z]$ denotes the weighted conditional linear predictor

$$\mathbb{E}_\Sigma^*[Y|X; Z] = X\mathbb{E}[X\Sigma(Z)^{-1}X'|Z]^{-1} \times \mathbb{E}[X\Sigma(Z)^{-1}Y|Z].$$

such misspecification the conditional moment restriction (6) will be invalid. Nevertheless the efficient score function may continue to have an expectation of zero at $\beta = \beta_0$. This suggest that an M-estimator based on an estimate of the efficient score function (7) may be consistent even if Assumption 1.5 does not hold. The following proposition provides conditions under which such a robustness property holds.

Proposition 3.1 (DOUBLE ROBUSTNESS) *Let $q_*(X) = q(X, \delta_*, h_*(X_2); \beta_0)$, $\rho_*(Z; \beta_0) = \psi(Z, \beta_0) - q_*(X)$ and redefine $H_0(X_2) = \mathbb{E} \left[\frac{\partial q_*(X)}{\partial h'} \middle| X_2 \right]$ and $G_0(X)$, $\Upsilon_0^h(X_2)$ and $\Sigma_0(X)$ similarly. Under restriction (1) and Assumptions 1.1 to 1.4 $\phi_\beta^f(Z, \eta, \beta_0)$ is mean zero if (i) $\eta = \eta_0$ and Assumption 1.5 holds or (ii) $\eta = \eta_* = (h_*, \delta_*, H_0, \Upsilon_0^h, G_0)$ and $p_0(x) = \Pr(D = 1 | X_2 = x_2) = e(x_2)$ for all $x \in \mathcal{X}$.*

Proof. Part (i) follows from Theorem 3.2. For part (ii) if $p_0(x) = e(x_2)$, then $G_0(X) = \frac{\partial q_*(X)}{\partial \delta'} - \mathbb{E}_\Sigma^* \left[\frac{\partial q_*(X)}{\partial \delta'} \middle| \frac{\partial q_*(X)}{\partial h'}; X_2 \right]$ and hence $\mathbb{E}[G_0(X)]$ is mean zero by construction. Furthermore

$$\begin{aligned} & \mathbb{E} \left[DH_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) \middle| X_2 \right] \\ &= \mathbb{E} \left[\mathbb{E}_\Sigma^* \left[\rho(Z; \beta_0) \middle| \frac{\partial q_0(X)}{\partial h'}; X_2 \right] \middle| X_2 \right] \\ &= \mathbb{E}[\psi(Z, \beta_0) | X_2] - \mathbb{E}[q(X, \delta_0, h_0(X_2); \beta_0) | X_2]. \end{aligned}$$

Therefore $\mathbb{E}[\phi_\beta^f(Z, \eta_*, \beta_0)] = \mathbb{E}[\psi(Z, \beta_0)] = 0$ as claimed. ■

Note that there is a tension between the robustness property of Proposition 3.1 and the size of the efficiency gain associated with Assumption 1.5. In particular when $p_0(x) = e(x_2)$, the second term in Ξ_0 will be identically equal to zero and there will be no benefits from parametric extrapolation (since the conditional distribution of X_1 given X_2 is the same in the $D = 1$ and $D = 0$ subpopulations). However it will still be the case that $\mathcal{I}_m^f(\beta_0) \geq \mathcal{I}_m(\beta_0)$ since, letting

$$\begin{aligned} R &= \frac{\rho(Z; \beta_0)}{e(X_2)^{1/2}} - \mathbb{E} \left[\frac{\partial q(X)}{\partial h'} \middle| X_2 \right] \\ &\quad \times \mathbb{E} \left[\left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right) \middle| X_2 \right]^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \frac{\rho(Z; \beta_0)}{e(X_2)^{1/2}}, \end{aligned}$$

we have $\mathbb{E}[RR'] = \mathbb{E}[\Lambda_0(X)] - \Xi_0$ equal to a positive semi-definite matrix. In this case Assumption 1.5 still delivers an efficiency gain from the dimension reduction in $q(X; \beta_0)$. If $X_1 = \emptyset$ and $X_2 = X$, such that $\mathbb{E}[\psi(Z, \beta_0) | x]$ is unrestricted, then $\mathcal{I}_m^f(\beta_0)$ simplifies to $\mathcal{I}_m(\beta_0)$ above. Therefore, Theorem 2.1 may be viewed as a special case of Theorem 3.1.

Appendices

A Proofs and derivations

To simplify notation in the Appendices let β denote the true parameter value β_0 unless explicitly stated otherwise (similarly the ‘0’ subscript is removed from other objects, such as the propensity score, when doing so does not cause confusion). All notation is as defined in the main text unless explicitly noted otherwise.

A.1 Proof of Theorem 2.1

The proof closely follows that of Theorem 1 in Chamberlain (1992b) and consists of three steps.

Step 1: Demonstration of equivalence with unconditional GMM problem The first step is to show that restrictions (4) and (5) are, in the multinomial case, equivalent to a finite set of unconditional moment restrictions. Under the multinomial assumption we have $X \in \{x_1, \dots, x_L\}$ for some L . Let the $L \times 1$ vector B have a 1 in the l^{th} row if $X = x_l$ and zeros elsewhere and $\tau_l = \Pr(X = x_l)$ (observe that $\sum_{l=1}^L \tau_l = 1$). Denote the value of the selection probability at $X = x_l$ by ρ_l and define $\rho = \{\rho_1, \dots, \rho_L\}'$; this vector gives the values of $p(\cdot)$ at each of the mass points of X . Using this notation we can write $p(X) = B'\rho$.

Under the multinomial assumption restrictions (4) and (5) are equivalent to the $L + K \times 1$ vector of unconditional moment restrictions

$$\mathbb{E}[m(Z, \beta, \rho)] = \mathbb{E} \begin{bmatrix} m_1(Z, \rho) \\ m_2(Z, \beta, \rho) \end{bmatrix} = \mathbb{E} \begin{bmatrix} B \left(\frac{D}{B'\rho} - 1 \right) \\ \frac{D}{B'\rho} \psi(Z, \beta) \end{bmatrix} = 0.$$

To verify that this is the case note that by iterated expectations

$$\mathbb{E}[m_1(Z, \rho)] = \begin{pmatrix} \tau_1 \mathbb{E} \left[\left(\frac{D}{p(X)} - 1 \right) \middle| X = x_1 \right] \\ \vdots \\ \tau_L \mathbb{E} \left[\left(\frac{D}{p(X)} - 1 \right) \middle| X = x_L \right] \end{pmatrix},$$

and hence $\mathbb{E}[m_1(Z, \rho)] = 0$ if and only if $\mathbb{E} \left[\frac{D}{p(X)} - 1 \middle| X \right] = 0$ for all $X \in \{x_1, \dots, x_L\}$. We also have

$$\mathbb{E}[m_2(Z, \beta, \rho)] = \mathbb{E} \left[\frac{D}{p(X)} \psi(Z, \beta) \right] = 0,$$

so $\mathbb{E}[m(Z, \beta, \rho)] = 0$ if and only if (4) and (5) are satisfied as claimed.

Step 2: Application of Lemma 2 of Chamberlain (1987) Chamberlain (1987, Lemma 2) shows that for Z a multinomial random variable the variance bound for β under the sole restriction that $\mathbb{E}[m(Z, \beta, \rho)] = 0$ is

$$\left\{ \left(M' V^{-1} M \right)^{-1} \right\}_{22}$$

where $\left\{ \left(M' V^{-1} M \right)^{-1} \right\}_{22}$ is the lower-right $K \times K$ block of $\left(M' V^{-1} M \right)^{-1}$ with

$$V \stackrel{def}{=} \mathbb{E} [m(Z, \beta, \rho) m(Z, \beta, \rho)'], \quad M \stackrel{def}{=} \mathbb{E} \left[\frac{\partial m(Z, \beta, \rho)}{\partial \rho'}, \frac{\partial m(Z, \beta, \rho)}{\partial \beta'} \right].$$

The application of Chamberlain's result requires that M has full column rank and that V is non-singular. The calculations made in Step 3 below demonstrate that these conditions are implied by the assumption that Γ has full column rank, $p(X)$ is bounded away from zero and non-singularity of $\mathbb{E}[\Lambda(X)]$.

Step 3: Calculation of the bound The final step is to solve for an explicit expression for $\left\{ \left(M' V^{-1} M \right)^{-1} \right\}_{22}$. This requires some simple, albeit tedious, algebra. Partitioning V_0

$$V_{L+K \times L+K} = \begin{pmatrix} V_{11} & V_{12} \\ V_{12}' & V_{22} \end{pmatrix}$$

we have the lower right-hand block, letting $\psi = \psi(Z, \beta)$ and $q(X) = \mathbb{E}[\psi|X]$, given by

$$\begin{aligned} V_{22}^{K \times K} &= \mathbb{E} [m_2(Z, \beta, \rho) m_2(Z, \beta, \rho)'] \\ &= \mathbb{E} \left[\frac{\mathbb{E} [\psi \psi' | X]}{p(X)} \right] \\ &= \mathbb{E} \left[\frac{\mathbb{V}(\psi|X)}{p(X)} + \frac{1-p(X)}{p(X)} q(X) q(X)' + q(X) q(X)' \right] \\ &= \sum_{l=1}^L \tau_l \left[\frac{\Sigma_l}{\rho_l} + \frac{1-\rho_l}{\rho_l} q_l q_l' + q_l q_l' \right], \end{aligned} \tag{8}$$

where $q_l = \mathbb{E}[\psi(Z, \beta) | x_l]$ and $\Sigma_l = \mathbb{V}(\psi | x_l)$.

The upper right-hand block is similarly derived as

$$\begin{aligned} V_{12}^{L \times K} &= \mathbb{E} [m_1(Z, \beta) m_2(Z, \beta, \rho)'] \\ &= \mathbb{E} \left[B \left(\frac{D}{B'\rho} - 1 \right) \left\{ \frac{D\psi(Z, \beta)}{B'\rho} \right\}' \right] \\ &= \mathbb{E} \left[B \left(\frac{1-p(X)}{p(X)} q(X)' \right) \right] \\ &= \left(\tau_1 \frac{1-\rho_1}{\rho_1} q_1 \quad \cdots \quad \tau_L \frac{1-\rho_L}{\rho_L} q_L \right)'. \end{aligned} \tag{9}$$

Finally the upper left-hand block is given by

$$\begin{aligned} V_{11}^{L \times L} &= \mathbb{E} \left[B \left(\frac{D}{B'\rho} - 1 \right) \left(\frac{D}{B'\rho} - 1 \right) B' \right] \\ &= \mathbb{E} \left[B B' \left(\frac{1-p(X)}{p(X)} \right) \right] \\ &= \text{diag} \left\{ \tau_1 \frac{1-\rho_1}{\rho_1} \quad \cdots \quad \tau_L \frac{1-\rho_L}{\rho_L} \right\}. \end{aligned} \tag{10}$$

Partition M

$$M_{L+K \times L+K} = \begin{pmatrix} M_{1\rho} & 0 \\ M_{2\rho} & M_{2\beta} \end{pmatrix},$$

where, from similar calculations to those made above, we have

$$M_{1\rho}^{L \times L} = -\text{diag} \left\{ \frac{\tau_1}{\rho_1} \quad \cdots \quad \frac{\tau_L}{\rho_L} \right\}, \quad M_{2\rho}^{K \times L} = - \left(\tau_1 \frac{q_1}{\rho_1} \quad \cdots \quad \tau_L \frac{q_L}{\rho_L} \right), \quad M_{2\beta}^{K \times K} = \Gamma. \tag{11}$$

Applying standard results on partitioned inverses then yields

$$M^{-1} = \begin{pmatrix} M_{1\rho}^{-1} & 0 \\ -M_{2\beta}^{-1}M_{2\rho}M_{1\rho}^{-1} & M_{2\beta}^{-1} \end{pmatrix},$$

Note that the existence of $M_{1\rho}^{-1}$ and $M_{2\beta}^{-1}$ follows from the assumptions that $p(X)$ is bounded away from zero and the assumption that Γ has full column rank.

Redundancy of knowledge of the propensity score suggests that $M^{-1}VM^{-1'}$ will be block diagonal. A sufficient condition for this is that (cf., Prokhorov and Schmidt 2006)

$$V'_{12} = M_{2\rho}M_{1\rho}^{-1}V_{11}. \quad (12)$$

To verify that this condition holds use (10) and (11) to show that

$$M_{2\rho}M_{1\rho}^{-1}V_{11} = \begin{pmatrix} \tau_1 \frac{1-\rho_1}{\rho_1} q_1 & \cdots & \tau_L \frac{1-\rho_L}{\rho_L} q_L \end{pmatrix},$$

which equals V'_{12} as required. Exploiting the resulting simplifications yields

$$M^{-1}VM^{-1'} = \begin{pmatrix} M_{1\rho}^{-1}V_{11}M_{1\rho}^{-1} & 0 \\ 0 & M_{2\beta}^{-1}(V_{22} - V'_{12}V_{11}^{-1}V_{12})M_{2\beta}^{-1'} \end{pmatrix}$$

and hence

$$(M^{-1}VM^{-1'})_{22} = M_{2\beta}^{-1}(V_{22} - V'_{12}V_{11}^{-1}V_{12})M_{2\beta}^{-1'}.$$

By $M_{2\rho}M_{1\rho}^{-1} = (q_1, \dots, q_L)$ and (12) we have $V'_{12}V_{11}^{-1}V_{12}$ equal to

$$\begin{aligned} V'_{12}V_{11}^{-1}V_{12} &= M_{2\rho}M_{1\rho}^{-1}V_{11}M_{1\rho}^{-1'}M'_{2\rho} \\ &= \sum_{l=1}^L \tau_l \frac{1-\rho_l}{\rho_l} q_l q_l' \\ &= \mathbb{E} \left[\frac{1-p(X)}{p(X)} q(X) q(X)' \right], \end{aligned}$$

and hence, using (8),

$$V_{22} - V'_{12}V_{11}^{-1}V_{12} = \mathbb{E} \left[\frac{\text{Var}(\psi|X)}{p(X)} + q(X) q(X)' \right] = \mathbb{E}[\Lambda(X)].$$

Using this result and taking the partitioned determinant gives

$$\det(V) = \det(V_{11}) \det(V_{22} - V'_{12}V_{11}^{-1}V_{12}) = \mathbb{E} \left[\frac{1-p(X)}{p(X)} \right] \det\{\mathbb{E}[\Lambda(X)]\},$$

and hence V is non-singular under overlap (Assumption 1.4) and non-singularity of $\mathbb{E}[\Lambda(X)]$.

Since $M_{2\beta} = \Gamma$ we have $\mathcal{I}_m(\beta_0) = \Gamma' \mathbb{E}[\Lambda(X)]^{-1} \Gamma$ as claimed. For completeness the upper left-hand portion of the full variance covariance matrix is given by

$$M_{11}^{-1}V_{11}M_{11}^{-1'} = \mathcal{I}_m^{-1}(\rho_0) = \text{diag} \left\{ \frac{p(x_1)(1-p(x_1))}{f(x_1)}, \dots, \frac{p(x_L)(1-p(x_L))}{f(x_L)} \right\}$$

where $f(x) = \sum_{l=1}^L \tau_l \times \mathbf{1}(x = x_l)$.

A.2 Proof of Theorem 3.1

The first two steps of the proof of Theorem 3.1 are analogous to those of Theorem 2.1 and therefore omitted. The actual calculation of the bound, while conceptually straightforward, is considerably more tedious. Details of this step are provided here.

Assume that the marginal distributions of X_1 and X_2 have I and M points of support with probabilities π_1, \dots, π_I and $\varsigma_1, \dots, \varsigma_M$. Let $L = I \times M$ and τ_{im} denote the joint probability $\Pr(X_1 = x_{1,i}, X_2 = x_{2,m})$. Let $\lambda = (\lambda_1, \dots, \lambda_M)'$ be the values of $h(\cdot)$ at each of the mass points of X_2 (for simplicity I assume that $\dim(h(x_2)) = P = 1$ in the calculations below, but the results generalize). Let C be a $M \times 1$ vector with a 1 in the m^{th} row if $X_2 = x_{2,m}$ and zeros elsewhere. Finally it is convenient to use the shorthand $\Psi = q(X)q(X)'$. In what follows I use both the single and double subscript notation to denote a point on the support of X as is convenient. We can map between the two notations by observing that $x_{im} = x_l$ for $l = (i-1)M + m$.

For the multinomial case the conditional moment problem defined by (4), (5) and (6) is equivalent to the unconditional problem

$$\mathbb{E}[m(Z, \theta)] = \mathbb{E} \begin{bmatrix} m_1(Z, \rho) \\ m_2(Z, \rho, \lambda, \delta, \beta) \\ m_3(Z, \rho, \beta) \end{bmatrix} = 0,$$

with $\theta = (\rho', \lambda', \delta', \beta')'$ and

$$\begin{aligned} m_1(Z, \rho) &= B \left(\frac{D}{B'\rho} - 1 \right), \quad m_2(Z, \rho, \lambda, \delta, \beta) = (B \otimes I_K) \left(\frac{D}{B'\rho} (\psi(Z, \beta) - q(X, \delta, C'\lambda; \beta)) \right), \\ m_3(Z, \rho, \beta) &= \frac{D}{B'\rho} \psi(Z, \beta). \end{aligned}$$

Partition $V = \mathbb{E}[m(Z, \theta)m(Z, \theta)']$ as

$${}_{L+KL+K \times L+KL+K}^V = \begin{pmatrix} V_{11} & & \\ V_{21} & V_{22} & \\ V_{31} & V_{32} & V_{34} \end{pmatrix},$$

where, using calculations similar to those given in the proof of Theorem 2.1, we have

$$\begin{aligned} V_{11} &= \text{diag} \left\{ \tau_1 \frac{1-\rho_1}{\rho_1}, \dots, \tau_L \frac{1-\rho_L}{\rho_L} \right\}, \quad V_{12} = (\underline{0}, \dots, \underline{0}), \\ V_{22} &= \text{diag} \left\{ \tau_1 \frac{\Sigma_1}{\rho_1}, \dots, \tau_L \frac{\Sigma_L}{\rho_L} \right\} \\ V_{31} &= \left(\tau_1 \frac{1-\rho_1}{\rho_1} q_1, \dots, \tau_L \frac{1-\rho_L}{\rho_L} q_L \right), \quad V_{32} = \left(\tau_1 \frac{\Sigma_1}{\rho_1}, \dots, \tau_L \frac{\Sigma_L}{\rho_L} \right) \\ V_{33} &= \sum_{l=1}^L \tau_l \left[\frac{\Sigma_l}{\rho_l} + \frac{1-\rho_l}{\rho_l} q_l q_l' + q_l q_l' \right]. \end{aligned}$$

We can partition the Jacobian matrix

$${}_{L+KL+K \times L+M+J+K}^M = \begin{pmatrix} M_{1\rho} & 0 & 0 & 0 \\ 0 & M_{2\lambda} & M_{2\delta} & 0 \\ M_{3\rho} & 0 & 0 & M_{3\beta} \end{pmatrix},$$

where

$$\begin{aligned}
M_{1\rho} &= -diag \left\{ \frac{\tau_1}{\rho_1}, \dots, \frac{\tau_L}{\rho_L} \right\} \\
M_{2\lambda} &= -(H'_1, \dots, H'_I)', \quad M_{2\delta} = - \begin{pmatrix} \tau_1 \nabla_\delta q_1 \\ \vdots \\ \tau_L \nabla_\delta q_L \end{pmatrix} \\
M_{3\rho} &= - \begin{pmatrix} \tau_1 \frac{q_1}{\rho_1} & \cdots & \tau_L \frac{q_L}{\rho_L} \end{pmatrix}, \quad M_{3\beta} = \Gamma.
\end{aligned}$$

where $H_i = diag \{ \tau_{i1} \nabla_h q_{i1}, \dots, \tau_{iM} \nabla_h q_{iM} \}$ for $i = 1, \dots, I$ with $q_{im} = q(x_{im}, \delta, h(x_{2,m}); \beta)$.

The variance bound for β is given by the lower right-hand $K \times K$ block of $(M'V^{-1}M)^{-1}$. We begin by calculating V^{-1} . Partition V

$$V = \begin{pmatrix} B_{11} & B_{12} \\ B'_{12} & B_{22} \end{pmatrix},$$

with

$$B_{11} = diag \{ V_{11} \quad V_{22} \}, \quad B_{12} = (V_{31} \quad V_{32})', \quad B_{22} = V_{33}.$$

Now partition V^{-1} as

$$V_0^{-1} = \begin{pmatrix} C_{11} & C_{12} \\ C'_{12} & C_{22} \end{pmatrix}, \quad (13)$$

where the partitioned inverse formula gives

$$C_{11} = diag \{ V_{11}^{-1} \quad V_{22}^{-1} \} + D' \mathbb{E}[\Psi]^{-1} D, \quad C'_{12} = -\mathbb{E}[\Psi]^{-1} D, \quad C_{22} = \mathbb{E}[\Psi]^{-1}$$

with $D = (A' \quad (\iota_L \otimes I_K)') = B'_{12} B_{11}^{-1}$ and $A = (q_1 \quad \cdots \quad q_L)'$ a $L \times K$ matrix.

Expression (13) follows since

$$\begin{aligned}
C_{22} &= (B_{22} - B'_{12} B_{11}^{-1} B_{12})^{-1} \\
&= \left\{ \sum_{l=1}^L \tau_l \left[\frac{\Sigma_l}{\rho_l} + \frac{1-\rho_l}{\rho_l} q_l q'_l + q_l q'_l \right] \right. \\
&\quad - \left(\tau_1 \frac{1-\rho_1}{\rho_1} q_1, \dots, \tau_L \frac{1-\rho_L}{\rho_L} q_L, \tau_1 \frac{\Sigma_1}{\rho_1}, \dots, \tau_L \frac{\Sigma_L}{\rho_L} \right) \\
&\quad \times diag \left\{ \tau_1 \frac{1-\rho_1}{\rho_1}, \dots, \tau_L \frac{1-\rho_L}{\rho_L}, \tau_1 \frac{\Sigma_1}{\rho_1}, \dots, \tau_L \frac{\Sigma_L}{\rho_L} \right\}^{-1} \\
&\quad \times \left(\tau_1 \frac{1-\rho_1}{\rho_1} q_1, \dots, \tau_L \frac{1-\rho_L}{\rho_L} q_L, \tau_1 \frac{\Sigma_1}{\rho_1}, \dots, \tau_L \frac{\Sigma_L}{\rho_L} \right)' \Big\}^{-1} \\
&= \sum_{l=1}^L \tau_l \left[\frac{\Sigma_l}{\rho_l} + \frac{1-\rho_l}{\rho_l} q_l q'_l + q_l q'_l \right] - \sum_{l=1}^L \tau_l \left[\frac{1-\rho_l}{\rho_l} q_l q'_l + \frac{\Sigma_l}{\rho_l} \right] \\
&= \left\{ \sum_{l=1}^L \tau_l q_l q'_l \right\}^{-1} = \mathbb{E}[\Psi]^{-1}.
\end{aligned}$$

We also have $C'_{12} = -C_{22} B'_{12} B_{11}^{-1} = -\mathbb{E}[\Psi]^{-1} D$ and

$$C_{11} = B_{11}^{-1} + B_{11}^{-1} B_{12} C_{22} B'_{12} B_{11}^{-1} = diag \{ V_{11}^{-1} \quad V_{22}^{-1} \} + D' \mathbb{E}[\Psi]^{-1} D.$$

We now evaluate $\mathcal{I}_m^f(\theta) = M'V^{-1}M$ to

$$\begin{pmatrix} M'_{1\rho}V_{11}^{-1}M_{1\rho} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ M'_{2\lambda} \begin{bmatrix} V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ -M'_{3\beta} \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \end{bmatrix} M_{2\lambda} & M'_{2\lambda} \begin{bmatrix} V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ -M'_{3\beta} \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \end{bmatrix} M_{2\delta} & -M'_{2\lambda} (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} M_{3\beta} & -M'_{2\delta} (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} M_{3\beta} \\ M'_{2\delta} \begin{bmatrix} V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ V_{22}^{-1} + (\iota_L \otimes I_K) \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \\ -M'_{3\beta} \mathbb{E}[\Psi]^{-1} (\iota_L \otimes I_K)' \end{bmatrix} M_{2\delta} & -M'_{3\beta} \mathbb{E}[\Psi]^{-1} M_{3\beta} & M'_{3\beta} \mathbb{E}[\Psi]^{-1} M_{3\beta} \end{pmatrix}$$

where I have made use of the equality $M'_{1\rho}A = M'_{3\rho}$.

Observe that, as in the standard semiparametric missing data model, $\mathcal{I}_m^f(\theta)$ satisfies Stein's condition for redundancy of knowledge of the propensity score for β . However the structure of the bound does indicate that knowledge of the finite dimensional parameters and nonparametric portions of the CEF of $\psi(Z, \beta)$ given X does increase the precision with which β can be estimated.

The variance bound for β_0 is given by the lower right-hand $K \times K$ block of the inverse of this matrix. Because of block diagonality we only need to consider the lower right-hand block. Partition this block as

$$\begin{pmatrix} B_{11} & B_{12} \\ B'_{12} & B_{22} \end{pmatrix}$$

where B_{11} , B_{12} and B_{22} are redefined to equal

$$\begin{aligned} B_{11} &= \begin{pmatrix} M'_{2\lambda}V_{22}^{-1}M_{2\lambda} & M'_{2\lambda}V_{22}^{-1}M_{2\delta} \\ M'_{2\delta}V_{22}^{-1}M_{2\lambda} & M'_{2\delta}V_{22}^{-1}M_{2\delta} \end{pmatrix} \\ &\quad + \begin{pmatrix} M'_{2\lambda}(\iota_L \otimes I_K) \\ M'_{2\delta}(\iota_L \otimes I_K) \end{pmatrix} \mathbb{E}[\Psi]^{-1} \begin{pmatrix} (\iota_L \otimes I_K)' M_{2\lambda} & (\iota_L \otimes I_K)' M_{2\delta} \end{pmatrix} \\ B_{12} &= \begin{pmatrix} M'_{2\lambda}(\iota_L \otimes I_K) \\ M'_{2\delta}(\iota_L \otimes I_K) \end{pmatrix} \mathbb{E}[\Psi]^{-1} M_{3\beta} \\ B_{33} &= M'_{3\beta} \mathbb{E}[\Psi]^{-1} M_{3\beta}. \end{aligned}$$

The information bound is therefore given by

$$\begin{aligned}
\mathcal{I}_m^f(\theta) &= B_{22} - B'_{12} B_{11}^{-1} B_{12} \\
&= M'_{3\beta} \mathbb{E}[\Psi]^{-1} M_{3\beta} - M'_{3\beta} \mathbb{E}[\Psi_0]^{-1} \\
&\quad \times \begin{pmatrix} (\iota_L \otimes I_K)' M_{2\lambda} & (\iota_L \otimes I_K)' M_{2\delta} \end{pmatrix} \\
&\quad \times \left\{ \begin{pmatrix} M'_{2\lambda} V_{22}^{-1} M_{2\lambda} & M'_{2\lambda} V_{22}^{-1} M_{2\delta} \\ M'_{2\delta} V_{22}^{-1} M_{2\lambda} & M'_{2\delta} V_{22}^{-1} M_{2\delta} \end{pmatrix} \right. \\
&\quad + \begin{pmatrix} M'_{2\lambda} (\iota_L \otimes I_K) \\ M'_{2\delta} (\iota_L \otimes I_K) \end{pmatrix} \mathbb{E}[\Psi]^{-1} \\
&\quad \times \begin{pmatrix} (\iota_L \otimes I_K)' M_{2\lambda} & (\iota_L \otimes I_K)' M_{2\delta} \end{pmatrix} \left. \right\}^{-1} \\
&\quad \times \begin{pmatrix} M'_{2\lambda} (\iota_L \otimes I_K) \\ M'_{2\delta} (\iota_L \otimes I_K) \end{pmatrix} \mathbb{E}[\Psi]^{-1} M_{3\beta} \\
&= M'_{3\beta} [\mathbb{E}[\Psi] + \begin{pmatrix} (\iota_L \otimes I_K)' M_{2\lambda} & (\iota_L \otimes I_K)' M_{2\delta} \end{pmatrix} \\
&\quad \times \begin{pmatrix} M'_{2\lambda} V_{22}^{-1} M_{2\lambda} & M'_{2\lambda} V_{22}^{-1} M_{2\delta} \\ M'_{2\delta} V_{22}^{-1} M_{2\lambda} & M'_{2\delta} V_{22}^{-1} M_{2\delta} \end{pmatrix}^{-1} \begin{pmatrix} M'_{2\lambda} (\iota_L \otimes I_K) \\ M'_{2\delta} (\iota_L \otimes I_K) \end{pmatrix}]^{-1} M_{3\beta},
\end{aligned}$$

where I have used the identity $A^{-1} - A^{-1}U(B^{-1} + U'A^{-1}U)^{-1}U'A^{-1} = (A + UBU')^{-1}$.

Using the partitioned inverse formula and multiplying out the expression in $[\cdot]$ above then gives

$$\begin{aligned}
\mathcal{I}_m^f(\theta) &= M'_{3\beta} \times [\mathbb{E}[\Psi] + (\iota_L \otimes I_K)' \left[M_{2\lambda} \left(M'_{2\lambda} V_{22}^{-1} M_{2\lambda} \right) M'_{2\lambda} \right. \\
&\quad + \left(M_{2\delta} - M_{2\lambda} \left(M'_{2\lambda} V_{22}^{-1} M_{2\lambda} \right)^{-1} M'_{2\lambda} V_{22}^{-1} M_{2\delta} \right) \\
&\quad \times \left(M'_{2\delta} V_{22}^{-1} M_{2\delta} - M'_{2\delta} V_{22}^{-1} M_{2\lambda} \left(M'_{2\lambda} V_{22}^{-1} M_{2\lambda} \right)^{-1} M'_{2\lambda} V_{22}^{-1} M_{2\delta} \right)^{-1} \\
&\quad \times \left. \left(M_{2\delta} - M_{2\lambda} \left(M'_{2\lambda} V_{22}^{-1} M_{2\lambda} \right)^{-1} M'_{2\lambda} V_{22}^{-1} M_{2\delta} \right) \right] (\iota_L \otimes I_K)]^{-1} \times M_{3\beta}.
\end{aligned}$$

We can now use the explicit expressions for V_0 and M_0 give above to generate an interpretable bound. The required calculations are tedious but straightforward (details are available in a supplemental Web Appendix), they give an information bound of $\mathcal{I}_m^f(\theta)$ as defined in the main text of the paper.

A.3 Proof of Theorem 3.2

In calculating the efficiency bound for the semiparametric missing data model defined by restriction (1) and Assumptions 1.1 to 1.5 above, I follow the general approach outlined by Bickel, Klaassen, Ritov and Wellner (1993) and, especially, Newey (1990, Section 3). First, I characterize the nuisance tangent space. Second, I demonstrate pathwise differentiability of the parameter of interest, β . The efficient influence function for this model equals the projection of the pathwise derivative onto the tangent space. In the present example the direct calculation this projection appears to be particularly difficult. However inspection of the variance bound associated with the conditional moment problem defined by restrictions (4), (5) and (6) provides a conjecture for the form of the efficient influence function. The third and final step of the proof therefore involves demonstrating that (i) this conjectured influence function lies in the model tangent space and (ii) that it is indeed the required projection (i.e., that it satisfies equation (9) in Newey (1990, p. 106)). The result then follows from an application of Theorem 3.1 in Newey (1990).

Step 1: Characterization of the nuisance tangent space Recalling that $Y = DY_1$, the joint density function for (Y, X, D) , making use of Assumption 1.3, is given by

$$f(y, x, d) = f(y_1 | x)^d p(x)^d [1 - p(x)]^{1-d} f(x).$$

Assumption 1.5 also requires that $f(y_1 | x)$ satisfy the restriction

$$\int \rho(z, \delta_0, h_0(x_2); \beta_0) f(y_1 | x) dy_1 = 0,$$

where

$$\rho(z, \delta, h(x_2); \beta) = \psi(y_1, x, \beta) - q(x, \delta, h(x_2); \beta).$$

Consider a regular parametric submodel with $f(y, x, d; \eta) = f(y, x, d)$ at $\eta = \eta_0$. The submodel joint density is given by

$$f(y, x, d; \eta) = f(y_1 | x; \eta)^d p(x; \eta)^d [1 - p(x; \eta)]^{1-d} f(x; \eta) \quad (14)$$

and satisfies the restriction

$$\int \rho(z, \delta(\eta), h(x_2; \eta); \beta_0) f(y_1 | x; \eta) dy_1 = 0. \quad (15)$$

The submodel score vector equals

$$s_\eta(y, x, d; \eta) = ds_\eta(y_1 | x; \eta) + \frac{d - p(x; \eta)}{p(x; \eta) [1 - p(x; \eta)]} \nabla_\eta p(x; \eta) + t_\eta(x; \eta), \quad (16)$$

where

$$\begin{aligned} s_\eta(y, x, d; \eta) &= \nabla_\eta \log f(y, x, d; \eta) \\ s_\eta(y_1 | x; \eta) &= \nabla_\eta \log f(y_1 | x; \eta), \quad t_\eta(x; \eta) = \nabla_\eta \log f(x; \eta). \end{aligned}$$

By the usual mean zero property of (conditional) scores we have

$$\mathbb{E}[s_\eta(Y_1 | X) | X] = \mathbb{E}[t_\eta(X)] = 0, \quad (17)$$

where suppression of η in a function means that it is evaluated at its population value (e.g., $t_\eta(x) = t_\eta(x; \eta_0)$).

Condition (15) imposes additional restrictions on $s_\eta(Y_1 | X)$ beyond conditional mean zeroness. To see the structure of these restrictions differentiate (15) with respect to η through the integral and evaluate the result at $\eta = \eta_0$:

$$\frac{\partial q_0(X)}{\partial \delta'} \frac{\partial \delta(\eta_0)}{\partial \eta'} + \frac{\partial q_0(X)}{\partial h'} \frac{\partial h(X_2; \eta_0)}{\partial \eta'} = \mathbb{E}[\rho(Z, \delta_0, h_0(X_2); \beta_0) s_\eta(Y_1 | X)' | X].$$

The conditional covariance between $\rho(Z, \delta_0, h_0(X_2); \beta_0)$ and $s_\eta(Y_1 | X)$ has a particular structure induced by the semiparametric restrictions on the form of $\mathbb{E}[\psi(Z, \beta) | x]$.

From (16), (17) and the above equality the tangent set is evidently

$$\mathcal{T} = \{ds(y_1 | x) + a(x)[d - p(x)] + t(x)\} \quad (18)$$

where $a(x)$ is unrestricted, $\mathbb{E}[s(Y_1 | X) | X] = \mathbb{E}[t(X)] = 0$ and

$$\mathbb{E}[\rho(Z, \delta_0, h_0(X_2); \beta_0) s(Y_1 | X)' | X] = \left(\frac{\partial q_0(X)}{\partial \delta'} \right) c + \left(\frac{\partial q_0(X)}{\partial h'} \right) k(X_2),$$

with c a constant matrix and $k(x_2)$ an unrestricted matrix-valued function of x_2 .

Step 2: Demonstration of pathwise differentiability Under the parametric submodel $\beta(\eta)$ is identified by the unconditional moment restriction

$$\mathbb{E}_\eta [\psi(Z; \beta(\eta))] = 0.$$

Differentiating under the integral and evaluating at $\eta = \eta_0$ gives

$$\frac{\partial \beta(\eta_0)}{\partial \eta'} = -\Gamma_0^{-1} \mathbb{E} \left[\psi(Z; \beta_0) \frac{\partial \log f(Y_1, X; \eta_0)}{\partial \eta'} \right]'$$

To demonstrate pathwise differentiability of β we require $F(Y, X, D)$ such that

$$\frac{\partial \beta(\eta_0)}{\partial \eta'} = \mathbb{E} [F(Y, X, D) s_\eta(Y, X, D)']$$

It is easy to verify that the function

$$F(Y, X, D) = -\Gamma_0^{-1} \left\{ \frac{D}{p_0(X)} \rho(Z, \delta_0, h_0(X_2); \beta_0) \right\} + q(X, \delta_0, h_0(X_2); \beta_0),$$

satisfies this condition (cf., Hahn 1998).

Step 3: Verification that conjectured efficient influence function equals the required projection Inspection of the variance bounds associated with the conditional moment problem suggests a candidate efficient influence function of

$$\begin{aligned} \phi_\beta^f(Z, \eta_0, \beta_0) &= \Gamma_0^{-1} \left\{ DH_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) \right. \\ &\quad \left. + D \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} G_0(X)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) + q(X; \beta_0) \right\}. \end{aligned} \quad (19)$$

Note that each term in the conjectured influence function is uncorrelated with the others.

I first verify that $\phi_\beta^f(Z, \eta_0, \beta_0)$ lies in the model tangent space. The last term in (19) plays the role of $t(x)$. A zero plays the role of $a(x)[d - p(x)]$. Finally the first two terms in (19) play the role of $ds(y|x)$. To see this note that in addition to being both conditionally mean zero we have

$$\begin{aligned} &\mathbb{E} \left[\rho(Z; \beta_0) \left\{ H_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) \right. \right. \\ &\quad \left. \left. + \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} G_0(X)' \Sigma_0(X)^{-1} \rho(Z; \beta_0) \right\}' \middle| X \right] \\ &= \left(\frac{\partial q_0(X)}{\partial \delta'} \right) c + \left(\frac{\partial q_0(X)}{\partial h'} \right) k(X_2) \end{aligned}$$

with

$$\begin{aligned} c &= \mathcal{I}_m^f(\delta_0)^{-1} \mathbb{E}[G_0(X)]' \\ k(X_2) &= \Upsilon_0^h(X_2)^{-1} \left\{ H_0(X_2)' - \Upsilon_0^{h\delta}(X_2) c \right\}. \end{aligned}$$

The candidate efficient influence function therefore belongs to the model tangent space as required.

I next show that $\phi_\beta^f(Z, \eta_0, \beta_0)$ is indeed the required projection by verifying that it satisfies

$$\mathbb{E} \left[\left\{ F(Y, X, D) - \phi_\beta^f(Z, \eta_0, \beta_0) \right\} t' \right] = 0, \text{ for all } t \in \mathcal{T}$$

(cf., equation (9) in Newey (1990, p. 106)). We have

$$F(Y, X, D) - \phi_\beta^f(Z, \eta_0, \beta_0) = \Gamma_0^{-1} D \left\{ \frac{1}{p_0(X)} - H_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \Sigma(X)^{-1} \right. \\ \left. - \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} G_0(X)' \Sigma(X)^{-1} \right\} \rho(Z; \beta_0).$$

By the conditional independence of Y_1 and D given X (Assumption 1.3) and conditional mean zeroness of $\rho(Z; \beta_0)$ it is easy to show that $F(Y, X, D) - \phi_\beta^f(Z, \eta_0, \beta_0)$ is orthogonal to any functions of the form $a(x)[d - p(x)]$ and $t(x)$. All that remains is to show orthogonality with $ds(y_1|x)$. We have

$$\mathbb{E} \left[\left\{ F(Y, X, D) - \phi_\beta^f(Z, \eta_0, \beta_0) \right\} Ds(Y_1|X)' \right] \\ = \mathbb{E} \left[\Gamma_0^{-1} \left\{ I_K - H_0(X_2) \Upsilon_0^h(X_2)^{-1} \left(\frac{\partial q_0(X)}{\partial h'} \right)' \left(\frac{\Sigma(X)}{p(X)} \right)^{-1} \right. \right. \\ \left. \left. - \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} G_0(X)' \left(\frac{\Sigma(X)}{p(X)} \right)^{-1} \right\} \right. \\ \left. \times \left\{ \left(\frac{\partial q_0(X)}{\partial \delta'} \right) c + \left(\frac{\partial q_0(X)}{\partial h'} \right) k(X_2) \right\} \right],$$

where I have made use of the special structure of the conditional covariance $\mathbb{E}[\rho(Z; \beta_0) s_\eta(Y_1|X)' | X]$. Multiplying out terms yields

$$\mathbb{E} \left[\left\{ F(Y, X, D) - \phi_\beta^f(Z, \eta_0, \beta_0) \right\} Ds(Y_1|X) \right] \\ = \Gamma_0^{-1} \mathbb{E} \left[\left\{ \frac{\partial q_0(X)}{\partial \delta'} c + H_0(X_2) k(X_2) \right. \right. \\ \left. \left. - H_0(X_2) \Upsilon_0^h(X_2)^{-1} \Upsilon_0^{h\delta}(X_2) c - H_0(X_2) k(X_2) \right. \right. \\ \left. \left. - \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} \Upsilon_0^\delta(X_2) c \right. \right. \\ \left. \left. + \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} \Upsilon_0^{h\delta}(X_2)' \Upsilon_0^h(X_2)^{-1} \Upsilon_0^{h\delta}(X_2) c \right. \right. \\ \left. \left. - \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} \Upsilon_0^{h\delta}(X_2)' k(X_2) \right. \right. \\ \left. \left. + \mathbb{E}[G_0(X)] \mathcal{I}_m^f(\delta_0)^{-1} \Upsilon_0^{h\delta}(X_2)' k(X_2) \right\} \right] \\ = \Gamma_0^{-1} \{ \mathbb{E}[G_0(X)] c - \mathbb{E}[G_0(X)] c \} = 0$$

where the first equality follows from iterated expectations and the second from the definitions of $G_0(X)$ and $\mathcal{I}_m^f(\delta_0)$ in the main text.

The result then follows from an application of Theorem 3.1 in Newey (1990).

References

- Angrist, Joshua D. and Alan B. Krueger. (1992). “The effect of age at school entry on educational attainment: an application of instrumental variables with moments from two samples,” *Journal of the American Statistical Association* 87 (418): 328 - 336.
- Bang, Heejung and James M. Robins. (2005). “Doubly robust estimation in missing data and causal inference models,” *Biometrics* 61 (4): 962 - 972.

- Bickel, Peter J., Chris A.J. Klaassen, Ya'acov Ritov and Jon A. Wellner. (1993). *Efficient and adaptive estimation for semiparametric models*. New York: Springer-Verlag, Inc.
- Brown, Bryan W. and Whitney K. Newey. (1998). "Efficient semiparametric estimation of expectations," *Econometrica* 66 (2): 453 - 464.
- Chamberlain, Gary. (1987). "Asymptotic efficiency in estimation with conditional moment restrictions," *Journal of Econometrics* 34 (1): 305 - 334.
- Chamberlain, Gary. (1992a). "Efficiency bounds for semiparametric regression," *Econometrica* 60 (3): 567 - 596.
- Chamberlain, Gary. (1992b). "Comment: sequential moment restrictions in panel data," *Journal of Business and Economic Statistics* 10 (1): 20 - 26.
- Chen, Xiaohong, Han Hong, Elie T. Tamer. (2005). "Measurement error models with auxiliary data," *Review of Economic Studies* 72 (2): 343 - 366.
- Chen, Xiaohong, Han Hong and Alessandro Tarozzi. (2004). "Semiparametric efficiency in GMM models of nonclassical measurement errors, missing data and treatment effects, *Mimeo*."
- Chen, Xiaohong, Han Hong and Alessandro Tarozzi. (2008). "Semiparametric efficiency in GMM models with auxiliary data," *Annals of Statistics* 36 (2): 808 - 843.
- Cheng, Philip E. (1994). "Nonparametric estimation of mean functionals with data missing at random," *Journal of the American Statistical Association* 89 (425): 81 - 87.
- Egel, Daniel, Bryan S. Graham, and Cristine Pinto. (2008). "Inverse probability tilting and missing data problems," *NBER Working Paper No. 13981*.
- Engle, Robert F., C. W. J. Granger, John Rice and Andrew Weiss. (1986). "Semiparametric estimates of the relation between weather and electricity sales," *Journal of the American Statistical Association* 81 (394): 310 - 320.
- Hahn, Jinyong. (1998). "On the role of the propensity score in efficient semiparametric estimation of average treatment effects," *Econometrica* 66 (2): 315 - 331.
- Hahn, Jinyong. (2004). "Functional restriction and efficiency in causal inference," *Review of Economics and Statistics* 86 (1): 73 - 76.
- Hájek, Jaroslav. (1972). "Local asymptotic minimax and admissibility in estimation," *Proceedings of the Sixth Berkeley Symposium on Mathematical Statistics and Probability* 1: 175 - 194 (L. M. Le Cam, J. Neyman & E. L. Scott, Eds.). Berkeley: University of California Press.
- Hirano, Keisuke, Guido W. Imbens and Geert Ridder. (2003). "Efficient estimation of average treatment effects using the estimated propensity score," *Econometrica* 71 (4): 1161 - 1189.

- Ichimura, Hidehiko and Oliver Linton. (2005). "Asymptotic expansions for some semiparametric program evaluation estimators," *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*: 149 -170 (D.W.K Andrews & J.H. Stock, Eds). Cambridge: Cambridge University Press.
- Imbens, Guido W. (2004). "Nonparametric estimation of average treatment effects under exogeneity: a review," *Review of Economics and Statistics* 86 (1): 4 - 29.
- Imbens, Guido W., Whitney K. Newey and Geert Ridder (2005). "Mean-square-error calculations for average treatment effects," *IEPR Working Paper 05.34*.
- Little, Roderick J. A. and Donald B. Rubin. (2002). *Statistical Analysis with Missing Data*. Hoboken, N.J.: John Wiley & Sons, Inc.
- Newey, Whitney K. (1990). "Semiparametric efficiency bounds," *Journal of Applied Econometrics* 5 (2): 99 - 135.
- Newey, Whitney K. (1994a). "Series estimation of regression functionals," *Econometric Theory* 10 (1): 1 - 28.
- Newey, Whitney K. (1994b). "The asymptotic variance of semiparametric estimators," *Econometrica* 62 (6): 1349 - 1382.
- Newey, Whitney K. (2004). "Efficient semiparametric estimation via moment restrictions," *Econometrica* 72 (6): 1877 - 1897.
- Newey, Whitney K. and Daniel McFadden. (1994). "Large sample estimation and hypothesis testing," *Handbook of Econometrics 4*: 2111 - 2245 (R.F. Engle & D.L. McFadden). Amsterdam: North Holland.
- Prokhorov, Artem and Peter J, Schmidt. (2006). "GMM redundancy results for general missing data problems," *Mimeo*.
- Robins, James M., Fushing Hsieh and Whitney Newey. (1995). "Semiparametric efficient estimation of a conditional density function with missing or mismeasured covariates," *Journal of the Royal Statistical Society B* 57 (2): 409 - 424.
- Robins, James M. and Andrea Rotnitzky. (1995). "Semiparametric efficiency in multivariate regression models," *Journal of the American Statistical Association* 90 (429): 122 - 129.
- Robins, James M., Andrea Rotnitzky and Lue Ping Zhao. (1994). "Estimation of regression coefficients when some regressors are not always observed," *Journal of the American Statistical Association* 89 (427): 846 - 866.
- Scharfstein, Daniel O., Andrea Rotnitzky and James M. Robins. (1999). "Rejoinder," *Journal of the American Statistical Association* 94 (448): 1135- 1146.
- Tsiatis, Anastasios A. (2006). *Semiparametric Theory and Missing Data*. New York: Springer.

Wang, Qihua, Oliver Linton and Wolfgang Härdle. (2004). “Semiparametric regression analysis with missing response at random,” *Journal of the American Statistical Association* 99 (466): 334 - 345.

Wooldridge, Jeffrey M. (1999a). “Asymptotic properties of weighted M-estimators for variable probability samples,” *Econometrica* 67 (6): 1385 - 1406.

Wooldridge, Jeffrey M. (1999b). “Distribution-free estimation of some nonlinear panel data models,” *Journal of Econometrics* 90 (1): 77 - 97.

Wooldridge, Jeffrey M. (2002). “Inverse probability weighted M-estimators for sample selection, attrition and stratification,” *Portuguese Economic Journal* 1 (2): 117 - 139.

Wooldridge, Jeffrey M. (2007). “Inverse probability weighted estimation for general missing data problems,” *Journal of Econometrics* 141 (2): 1281 - 1301.