

Semiparametric Efficiency Bounds for Conditional Moment Restriction Models with Different Conditioning Variables

Marian Hristache & Valentin Patilea
CREST (Ensaï)*

October 16, 2014

Abstract

This paper addresses the problem of semiparametric efficiency bounds for conditional moment restriction models with different conditioning variables. We characterize such an efficiency bound, that in general is not explicit, as a limit of explicit efficiency bounds for a decreasing sequence of unconditional (marginal) moment restriction models. An iterative procedure for approximating the efficient score when this is not explicit is provided. Our theoretical results provide new insight for the theory of semiparametric efficiency bounds literature and open the door to new applications. In particular, we investigate a class of regression-like (mean regression, quantile regression,...) models with missing data, an example of demand-supply simultaneous equations model and a generalized bivariate dichotomous model.

KEY WORDS : Conditional moment restrictions models; Semiparametric efficiency bounds; Tangent spaces; Missing data models; Simultaneous equations models.

1 The model

Conditional moment restrictions models represent a large class of statistical models. Seemingly unrelated nonlinear regressions, see Gallant (1975), Müller (2009), seemingly unrelated quantile regressions, see Jun and Pinske (2009), regression models with missing data, see Robins, Rotnitzky and Zhao (1994), Tsiatis (2006), are

*CREST, Ecole Nationale de la Statistique et de l'Analyse de l'Information (Ensaï), Campus de Ker-Lann, rue Blaise Pascal, BP 37203, 35172 Bruz cedex, France.

Authors emails: marian.hristache@ensai.fr, valentin.patilea@ensai.fr

only few examples and related contributions. Many other references and examples of statistical and econometric models that could be stated as conditional moment restrictions models are provided in Ai and Chen (2012) and Hansen (2008).

In this paper we address the problem of calculating semiparametric efficiency bounds in models defined by several conditional moment restrictions with possibly different conditioning variables. More formally, the sample under study consists of independent copies of a random vector $Z \in \mathcal{Z} \subset \mathbb{R}^q$. Let J be some positive integer that is fixed in the following. For any $j \in \{1, \dots, J\}$, let $X^{(j)}$ be a random q_j -dimension subvector of Z , where $0 \leq q_j < q$. Let $g_j : \mathbb{R}^q \times \mathbb{R}^d \rightarrow \mathbb{R}^{p_j}$, $j \in \{1, \dots, J\}$, denote given functions of Z and the unknown parameter $\theta \in \Theta \subset \mathbb{R}^d$. The semiparametric model we consider is defined by the conditional moment restrictions

$$E \left[g_j(Z, \theta) \mid X^{(j)} \right] = 0, \quad j = 1, \dots, J, \quad \text{almost surely (a.s.).} \quad (1)$$

Such models are common in the econometric literature, see for instance the equations (14.33) and (14.36) in Wooldridge (2010). As usually, the notation $g_j(Z, \theta)$ does not necessarily mean that the function g_j depends on all the components of Z . It is assumed that the d -dimension parameter θ is identified by the conditional restrictions, which means there exists a unique value θ_0 such that the true law of Z satisfies equations (1). By definition, $X^{(j)}$ is a constant random variable when $q_j = 0$, and hence the conditional expectation given $X^{(j)}$ is the marginal expectation.

Particular cases of this model have been extensively studied in the literature. For $J = 1$ and $q_1 = 0$ we obtain a model defined by an unconditional set of moment equations

$$E[g(Z, \theta)] = 0.$$

Hansen (1982) considered the class of GMM estimators and showed how to construct an optimal one in that class. Its asymptotic variance equals the the semiparametric efficiency bound obtained by Chamberlain (1987).

The GMM method extends naturally to models defined by conditional moment equations, corresponding to the case $J = 1$ and $q_1 > 0$ in our setting, that is

$$E[g(Z, \theta) \mid X] = 0 \quad (\text{a.s.}).$$

From a mathematical point of view, such a model is equivalent to the intersection of the models of the form

$$E[a(X) g(Z, \theta)] = 0,$$

where $a(X)$ is an arbitrary conformable random matrix whose entries are square integrable. In the econometric literature $a(X)$ is referred to as a matrix of *instruments*. The supremum of the information on θ_0 in these models yields the semiparametric Fisher information on θ_0 in the conditional equation model, obtained by

Chamberlain (1992a). It is also the information on θ_0 for the unconditional moment equation

$$E [a^* (X) g (Z, \theta)] = 0,$$

with properly chosen ‘optimal’ instruments $a^* (X)$.

A further generalization, which can also be written under the form (1), is given by a sequential (nested) moment restrictions model, in which the σ -fields generated by the conditioning vectors satisfy the condition $\sigma (X^{(1)}) \subset \sigma (X^{(2)}) \subset \dots \subset \sigma (X^{(J)})$. For the expression of the semiparametric efficiency bound in the sequential case, see Chamberlain (1992b) and Ai and Chen (2012); see also Hahn (1997) and Ahn and Schmidt (1999) and references therein for examples of applications. It turns out that once again the information on θ_0 can be obtained by taking the supremum of the information on θ_0 in the following unconditional models :

$$E \left[a_j \left(X^{(j)} \right) g_j (Z, \theta) \right] = 0, \quad j = 1, \dots, J,$$

where the number of lines of the matrices a_j is fixed and equal to the dimension of θ and the supremum is attained for a suitable choice $a_1^* (X^{(1)}), \dots, a_J^* (X^{(J)})$ of optimal instruments. The reason why this happens in the case with nested σ -fields is the fact that the model of interest can be written as the decreasing limit of a sequence of models for which a so-called ‘*spanning condition*’, similar to the one considered in Newey (2004), holds and the limit of the corresponding efficient scores has an explicit solution.

In this paper we show that the information on θ_0 in model (1) can be obtained as the limit of the information on θ_0 in a decreasing sequence of unconditional moment models of the form

$$E \left[a_j^{(k)} \left(X^{(j)} \right) g_j (Z, \theta) \right] = 0, \quad j = 1, \dots, J, \quad k = 1, 2, \dots, \quad (2)$$

where the numbers of lines in the matrices $a_j^{(k)}$ increases to infinity with k . To our best knowledge this result is new. It provides theoretical support for a natural solution that could be used in practice: replace the model (1) by a large number of unconditional moment conditions like in equation (2) in order to approach efficiency. Herein we also propose an alternative route for approximating the efficiency bound. More precisely, we give a general method to approximate the efficient score, which in most of the situations does not have an explicit form as in the aforementioned examples. In particular, our general approach for approximating the efficient score brings in a new light on the functional equations used to characterize the efficient score in the regression model with unobserved explanatory variables in Robins, Rotnitzky and Zhao (1994); see also Tsiatis (2006) and Tan (2011).

To summarize, our theoretical results provide new insight for the theory of semi-parametric efficiency bounds literature and open the door to new applications, in particular for nonlinear simultaneous equations models and in missing data contexts. Quantile regression models with missing covariables is one example of such new application that could be treated in our theoretical framework. This case of practical interest is different and more difficult to handle than the quantile regression with missing responses as considered by Wei, Ma and Carroll (2012); see also Müller and van Keilegom (2012). On the other hand, our methodology allows the investigation of nonlinear simultaneous equations models in greater generality.

The paper is organized as follows. Section 2 contains our main results. We show that under a suitable ‘spanning condition’ on the tangent spaces, the semiparametric Fisher information in model (1) can be obtained as the limit of the efficiency bounds for a decreasing sequence of models. In section 3 we propose a ‘backfitting’ procedure, for computing the projection of the score on the tangent space of the model. With at hand an approximation of the efficient score, we suggest a general method for constructing asymptotically efficient estimators. In section 4 we illustrate the utility of our theoretical results for four classes of models: sequential (nested) conditional models, regression-like models with missing data, simultaneous equations models and generalized dichotomous models. Some technical results and proofs are relegated to the Appendix.

2 The main results

Let us introduce some notation and definitions, see also van der Vaart (1998), sections 25.2 and 25.3. Given a sample space \mathcal{Z} and a probability P on the sample space, we denote by $L^2(P)$ the usual Hilbert space of measurable real-valued functions that are squared-integrable with respect to P . For \mathcal{H} a Hilbert space and $\mathcal{S} \subset \mathcal{H}$ let $\bar{\mathcal{S}}$ denote the closure of \mathcal{S} in \mathcal{H} . The statistical models on the sample space \mathcal{Z} , are denoted by $\mathcal{P}, \mathcal{P}_1, \mathcal{P}_2, \dots$. A statistical model is a collection of probability measures defined by their densities with respect to some fixed dominating measure on the sample space. Herein the vectors are column matrices and $A \in \mathbb{R}^r \times \mathbb{R}^s$ means A is a $r \times s$ -matrix with random elements, if not stated differently. $E(A)$ denotes the expectation of A . For $A \in \mathbb{R}^r \times \mathbb{R}^r$, $E^{-1}(A)$ denotes the inverse of the square matrix $E(A)$, whenever the inverse exists. Finally, for a square matrix A , let A^- denote a generalized inverse, for instance the Moore-Penrose pseudoinverse.

2.1 Efficiency bound

The main idea we follow to derive the semiparametric efficiency bound for the parameter θ_0 is to transform the finite number of conditional moment restrictions (1)

in a countable number of unconditional (marginal) moment restrictions. Next, for any finite subset of these unconditional moment restrictions, one could easily obtain the Fisher information bound. Eventually, one may expect to obtain the semiparametric efficiency bound for the model (1) as the limit of the efficiency bounds for a decreasing sequence of models defined by an increasing sequence of finite subsets of unconditional moment restrictions. Remark 1 in the Appendix proves that in general this intuition is not correct. However, the subsequent Lemma 2 states that this intuition becomes correct under the additional ‘spanning condition’ formally stated in equation (26) in the Appendix.

Let us introduce some more notation. For a vector a , we denote by a' its transpose. If $\zeta : \mathcal{Z} \times \Theta \rightarrow \mathbb{R}^m$, $m \geq 1$, is some given function of Z and θ and X is some subvector of Z , we denote

$$\partial_{\theta'} E[\zeta | X] = \partial_{\theta'} E[\zeta(Z, \theta_0) | X] = \left. \frac{\partial}{\partial \theta'} E[\zeta(Z, \theta) | X] \right|_{\theta=\theta_0} \in \mathbb{R}^d \times \mathbb{R}^m, \quad (3)$$

when such derivatives of the map $\theta \mapsto E[\zeta(Z, \theta) | X]$ exist. A similar notation will be used with the conditional expectation $E(\cdot | X)$ replaced by the marginal (unconditional) expectation with respect to the law of Z . Let us point out that the maps $\theta \mapsto \zeta(z, \theta)$ may not be everywhere differentiable. Next, let us define

$$\underline{g} = (g'_1, \dots, g'_J)' \in \mathbb{R}^p, \quad \text{where } p = p_1 + \dots + p_J,$$

and let \underline{X} denote the vector of all components of Z contained in the subvectors $X^{(j)}$, $j = 1, \dots, J$. Let P_0 be the true law of the vector Z and $P_{X^{(j)}}$ the law of $X^{(j)}$, $j = 1, \dots, J$.

For the purpose of transforming conditional moments in unconditional versions, let us consider J countable sets of squared integrable functions $\mathcal{W}^{(j)} = \{w_k^{(j)} : k \in \mathbb{N}^*\} \subset L^2(P_{X^{(j)}})$ such that $\overline{\text{lin}} \mathcal{W}^{(j)} = L^2(P_{X^{(j)}})$, that is the linear span of $\mathcal{W}^{(j)}$ is dense in $L^2(P_{X^{(j)}})$, for $1 \leq j \leq J$. For any $k \in \mathbb{N}^*$, let

$$\underline{g}_k^w(Z, \theta) = \begin{pmatrix} w_{(k)}^{(1)}(X^{(1)}) \otimes g_1(Z, \theta) \\ w_{(k)}^{(2)}(X^{(2)}) \otimes g_2(Z, \theta) \\ \vdots \\ w_{(k)}^{(J)}(X^{(J)}) \otimes g_J(Z, \theta) \end{pmatrix},$$

where \otimes denotes the Kronecker product and

$$w_{(k)}^{(j)}(X^{(j)}) = (w_1^{(j)}(X^{(j)}), \dots, w_k^{(j)}(X^{(j)}))' \in \mathbb{R}^k.$$

Let $I_{\theta_0}^{(k)}$ be the Fisher information on θ_0 in the model

$$E \left[\underline{g}_k^w (Z, \theta) \right] = 0, \quad (4)$$

that is

$$I_{\theta_0}^{(k)} = \left[\partial_{\theta'} E \left(\underline{g}_k^w (Z, \theta_0) \right) \right]' V^{-1} \left[\underline{g}_k^w (Z, \theta_0) \right] \partial_{\theta'} E \left(\underline{g}_k^w (Z, \theta_0) \right).$$

See Chamberlain (1987), Newey (2001), see also Chen and Pouzo (2009) for the non-smooth case.

We can state now the main result of the paper.

Theorem 1 *Under the Assumptions T and SP in section 2.2, the information bound I_{θ_0} on θ_0 at P_0 in model (1) is given by*

$$I_{\theta_0} = \lim_{k \rightarrow \infty} I_{\theta_0}^{(k)},$$

where, for any $k \in \mathbb{N}^*$, $I_{\theta_0}^{(k)}$ is the Fisher information on θ_0 in the model defined as in (4).

In the general theory of efficiency bounds, the semiparametric Fisher information on a finite dimension parameter in a semiparametric model is the infimum of the Fisher information over all its parametric submodels; see for instance Newey (1990). For models defined by conditional moment equations, Theorem 1 shows that the same semiparametric Fisher information can be alternatively obtained as the lower limit of the semiparametric Fisher information in a sequence of decreasing supra-models. The main reason for this is that with such a decreasing sequence of supra-models, the ‘spanning condition’ (26) holds true. Moreover, since $L^2(P_0)$ is a separable Hilbert space, Theorem 1 can be restated under the following equivalent form.

Let

$$\mathcal{B} = \left\{ \left(\omega_1 \left(X^{(1)} \right), \dots, \omega_J \left(X^{(J)} \right) \right) : \omega_{j,lk} \in L^2 \left(P_{X^{(j)}} \right), \right. \\ \left. 1 \leq l \leq d, 1 \leq k \leq p_j, 1 \leq j \leq J \right\};$$

that is, for any j , $\omega_j(X^{(j)})$ is a $d \times p_j$ matrix of functions whose entries $\omega_{j,lk}(X^{(j)})$ are squared integrable, for $1 \leq l \leq d$ and $1 \leq k \leq p_j$. For a $d \times p$ -matrix $\omega = \omega(\underline{X}) = (\omega_1(X^{(1)}), \dots, \omega_J(X^{(J)})) \in \mathcal{B}$, let $I_{\theta_0}(\omega)$ be the Fisher information on θ_0 in the model defined by the marginal moment restrictions

$$E \left[\omega_j \left(X^{(j)} \right) g_j (Z, \theta) \right] = 0, \quad j \in \{1, \dots, J\}, \quad (5)$$

model which can also be written under the compact form $E \left[\omega(\underline{X}) \underline{g}(Z, \theta) \right] = 0$.

Corollary 1 *Under the conditions of Theorem 1, we have*

$$I_{\theta_0} = \sup_{\omega \in \mathcal{B}} I_{\theta_0}(\omega).$$

The proof of Theorem 1 is based on the fact that the tangent space for the model (1) is the intersection of the tangent spaces corresponding to the J conditional moment models defined by each of the vector functions g_1, \dots, g_J . (See the Appendix for the formal definition of a tangent space.) This result is of own interest and hence we state it below. Moreover, it will be used in section 3 to propose an alternative strategy for approximating the efficient score.

Proposition 1 *The tangent space (of the nonparametric part) of the model (1) defined by $E[g_j(Z, \theta) | X^{(j)}] = 0$ (a.s.), $\forall j = 1, \dots, J$, is the intersection of the tangent spaces (of the nonparametric parts) of the models defined, for each $j \in \{1, \dots, J\}$, by*

$$E[g_j(Z, \theta) | X^{(j)}] = 0 \text{ (a.s.)}.$$

2.2 Assumptions

Assumption T There exist bounded squared integrable vector functions b_1, \dots, b_J (with the same dimensions as g_1, \dots, g_J , respectively) such that:

1. for any $i, j \in \{1, \dots, J\}$ with $i \neq j$, $E(g_i(Z, \theta_0) b'_j(Z) | X^{(i)}, X^{(j)}) = 0$ (a.s.);
2. for any $i \in \{1, \dots, J\}$, $E(b_i(Z) g'_i(Z, \theta_0) | X^{(i)})$ is invertible (a.s.) and

$$c_b = \max_{1 \leq i \leq J} \left\| E^{-1} \left(b_i(Z) g'_i(Z, \theta_0) | X^{(i)} \right) \right\|_{\infty} < \infty;$$

3. $\max_{1 \leq i \leq J} \left\| E \left(g'_i(Z, \theta_0) g_i(Z, \theta_0) | X^{(i)} \right) E(\mathbf{1}\{b_i(Z) \neq 0\} | X^{(i)}) \right\|_{\infty} < \infty$ (a.s.).

Assumption SP 1. The models \mathcal{P} defined by (1) and \mathcal{P}_k defined by (4), with $k \in \mathbb{N}^*$, can be written in the semiparametric form

$$\mathcal{P} = \{P_{\theta, \eta} : \theta \in \Theta, \eta \in H\}, \quad \mathcal{P}_k = \{P_{\theta, \eta} : \theta \in \Theta, \eta \in H_k\}, \quad k \in \mathbb{N}^*,$$

and satisfy the assumptions of Lemma 25.25 of van der Vaart (1998).

2. The Fisher information matrices I_{θ_0} and $I_{\theta_0}^{(k)}$ on θ_0 in models \mathcal{P} and \mathcal{P}_k respectively, for any $k \in \mathbb{N}^*$, are well defined and nonsingular.

Primitive conditions guaranteeing Assumption T are provided in the technical Assumption T' in the Appendix. Let us comment on the role of Assumption T. It ensures that $\mathcal{M} \cap \mathcal{T}'$ is dense in \mathcal{T}' where \mathcal{M} is the subspace of bounded functions of Z . This density condition is commonly used in the literature to show that the candidate tangent space \mathcal{T}' of the nonparametric part of the model is indeed that tangent space. However, on contrary of what is sometimes implicitly supposed in the literature, there is no general mathematical result that guarantees that the subspace of bounded functions is dense in any subspace of the squared integrable functions of Z .¹ Our proofs also rely on the fact that $\mathcal{M} \cap \mathcal{T}'$ is dense in \mathcal{T}' and Assumption T provides a general sufficient condition guaranteeing this property with $J \geq 1$. In the case of bounded and orthogonal g_j 's, one could simply take $b_j = g_j$, $1 \leq j \leq J$.

To guarantee Assumption SP.2 it suffices to suppose that for any $1 \leq j \leq J$: (i) $\|V(g_j(Z, \theta_0) | X^{(j)})\|_\infty < \infty$; (ii) the maps $\theta \mapsto E(g_j(Z, \theta_0) | X^{(j)} = x^{(j)})$ are differentiable for $P_{X^{(j)}}$ -almost all $x^{(j)}$; and (iii) the information matrix

$$E \left\{ \left[\partial_{\theta'} E \left(g_j(Z, \theta_0) | X^{(j)} \right) \right]' V^{-1} \left[g_j(Z, \theta_0) | X^{(j)} \right] \partial_{\theta'} E \left(g_j(Z, \theta_0) | X^{(j)} \right) \right\}$$

is non singular.

A consequence of Assumption SP (see Lemma 25.25 of van der Vaart (1998)) is that the parameter defined by $\psi(P_{\theta, \eta}) = \theta$ is differentiable at $P_0 = P_{\theta_0, \eta_0}$ with respect to the tangent space $\mathcal{T} = \mathcal{T}(\mathcal{P}, P_0)$. It also ensures that the tangent space \mathcal{T} can be written as the sum of the finite dimensional subspace spanned by the components of the parametric score S_{θ_0} and the tangent space \mathcal{T}' corresponding to the nonparametric part $\mathcal{P}' = \{P_{\theta_0, \eta} : \eta \in H\}$ of the model \mathcal{P} :

$$\mathcal{T} = \text{lin}S_{\theta_0} + \mathcal{T}'.$$

3 Approximating the efficient score

To simplify the presentation, in this section let us consider only the case $J = 2$. To obtain an efficient estimator, a common way is to solve θ from the efficient score

¹Here we provide a counterexample showing that the set of bounded functions is not dense in any subspace of $L^2(P_0)$. Indeed, if $\text{supp}P_0 = \text{supp}Z = [0, 1]^q \subset \mathbb{R}^q$, let $\phi_0(u) = \|u\|^{-1/4} \mathbf{1}\{0 < \|u\| < 1\}$ and take \mathcal{S} the linear span of $\{\phi_0(\|Z\| - k), k \in \mathbb{Z}\}$. Then \mathcal{S} is an infinite dimensional subspace of $L^2(P_0)$ and the subspace $\mathcal{M} \subset L^2(P_0)$ of bounded functions is not dense in \mathcal{S} since $\mathcal{M} \cap \mathcal{S} = \{0\}$, so that $\overline{\mathcal{M} \cap \mathcal{S}} = \{0\}$.

equations; see van der Vaart (1998), section 25.8. By definition, the efficient score is the componentwise projection of the score S_{θ_0} on the orthogonal complement of the tangent space $\mathcal{T} = \mathcal{T}(\mathcal{P}, P_0)$ defined in equation (35). In the projection of S_{θ_0} on \mathcal{T}^\perp only the nonparametric part of the tangent space matters. Moreover, the projection of S_{θ_0} is componentwise. It is then common practice in the literature to identify \mathcal{T} with the subspace of $\{L^2(P_0)\}^d = \bigoplus_{k=1}^d L^2(P_0)$ obtained as the d -fold cartesian product of the nonparametric part of \mathcal{T} . Here the direct sum of Hilbert spaces is considered with the usual inner product $\langle (\phi_1, \dots, \phi_d), (\psi_1, \dots, \psi_d) \rangle = \langle \phi_1, \psi_1 \rangle + \dots + \langle \phi_d, \psi_d \rangle$. Therefore we will slightly change our notation for the tangent spaces. More precisely, let us define

$$\begin{aligned} \mathcal{T} &= \left\{ s \in \bigoplus_{k=1}^d L^2(P_0) : E(s) = 0, E\left(g_i(Z, \theta_0) s'(Z) \mid X^{(i)}\right) = 0, i = 1, 2 \right\} \\ &= \mathcal{T}_1 \cap \mathcal{T}_2, \end{aligned}$$

where, for $i = 1, 2$,

$$\mathcal{T}_i = \left\{ s \in \bigoplus_{k=1}^d L^2(P_0) : E(s) = 0, E\left(g_i(Z, \theta_0) s'(Z) \mid X^{(i)}\right) = 0 \right\},$$

so that

$$\mathcal{T}_i^\perp = \left\{ s \in \bigoplus_{k=1}^d L^2(P_0) : s(Z) = a_i\left(X^{(i)}\right) g_i(Z, \theta_0) \right\},$$

where a_i is an arbitrary matrix-valued function (of compatible dimension $d \times p_i$) whose entries are squared integrable functions of $X^{(i)}$, $i = 1, 2$. Clearly, $\mathcal{T}^\perp = \overline{\mathcal{T}_1^\perp + \mathcal{T}_2^\perp}$.

In general, the projection of S_{θ_0} on \mathcal{T}^\perp is not explicit. To approximate this projection and to further build an asymptotically efficient estimator for model (1), two strategies can be developed from the results obtained in this paper. The first one is based on Theorem 1. With the same notations as in section 2.1, let \mathcal{P}_k be the model defined by (4), that is

$$E\left[w_{(k)}^{(j)}\left(X^{(j)}\right) \otimes g_j(Z, \theta)\right] = 0, \quad j = 1, \dots, J.$$

Then $\{\mathcal{P}_k\}_{k \in \mathbb{N}^*}$ is a decreasing sequence of models and the corresponding tangent spaces satisfy (see the proof of Theorem 1)

$$\mathcal{T}_1 \supset \mathcal{T}_2 \supset \dots \supset \mathcal{T}_k \supset \mathcal{T}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{T}_k = \mathcal{T},$$

where \mathcal{T} is the tangent space of model (1). Taking the sequence of the efficient scores $\bar{S}_{\theta_0}^{(k)}$ in models \mathcal{P}_k , Theorem 4.5 from Hansen and Sargent (1991) ensures the convergence (in $L^2(P_0)$) of $\bar{S}_{\theta_0}^{(k)}$ to the efficient score \bar{S}_{θ_0} of model (1). In conclusion, using sufficiently many suitable instruments in the original model (1) allows to obtain an approximate score $\bar{S}_{\theta_0}^{(k)}$ arbitrarily close to \bar{S}_{θ_0} . Results of this type have already been discussed in the literature in some particular cases of model (1); see, among others, Chamberlain (1992b), Newey (1993), Hahn (1997), Han and Phillips (2006). How many such instruments to choose, in particular how to define the dependence of k on the sample size n , is an open question in our general setting and will be a subject for future work.

Here, we will investigate an alternative approach, based on the form of the tangent space given in Proposition 1, or, to be more precise, based on the aforementioned writing of the orthogonal tangent space $\mathcal{T}^\perp = \overline{\mathcal{T}_1^\perp + \mathcal{T}_2^\perp}$. For this purpose, we will use the iterative ('backfitting' or successive approximation) procedure considered in Theorem A.4.2 of Bickel, Klaassen, Ritov and Wellner (1993), page 438; BKRW hereafter. Let $H_i = \mathcal{T}_i^\perp$, $g_i = g(Z, \theta_0)$, $i = 1, 2$, and let $\partial_\theta E(g'_i)$ be the transposed of the matrix $\partial_\theta E(g_i)$ defined in equation (3). The steps of the procedure we propose are the following :

1. Set $m = 0$. Take $a_1^{(0)} = 0$.
2. Put $m = m + 1$. Calculate

$$\bar{S}_{\theta_0}^{(m)} = a_1^{(m)} \left(X^{(1)} \right) g_1 + a_2^{(m)} \left(X^{(2)} \right) g_2$$

where

$$\begin{aligned} a_1^{(m)} \left(X^{(1)} \right) &= a_1^{(m)} \left(X^{(1)}, \theta_0 \right) = -\partial_\theta E \left(g'_1 \mid X^{(1)} \right) V^- \left(g_1 \mid X^{(1)} \right) \\ &+ E \left[\partial_\theta E \left(g'_2 \mid X^{(2)} \right) V^- \left(g_2 \mid X^{(2)} \right) g_2 g'_1 \mid X^{(1)} \right] V^- \left(g_1 \mid X^{(1)} \right) \\ &+ E \left[E \left[a_1^{(m-1)} \left(X^{(1)} \right) g_1 g'_2 \mid X^{(2)} \right] V^- \left(g_2 \mid X^{(2)} \right) g_2 g'_1 \mid X^{(1)} \right] V^- \left(g_1 \mid X^{(1)} \right) \end{aligned}$$

and

$$\begin{aligned} a_2^{(m)} \left(X^{(2)} \right) &= a_2^{(m)} \left(X^{(2)}, \theta_0 \right) = -\partial_\theta E \left(g'_2 \mid X^{(2)} \right) V^- \left(g_2 \mid X^{(2)} \right) \\ &- E \left[a_1^{(m)} \left(X^{(1)} \right) g_1 g'_2 \mid X^{(2)} \right] V^- \left(g_2 \mid X^{(2)} \right). \end{aligned}$$

3. Repeat from step 2 till the convergence of $\bar{S}_{\theta_0}^{(m)}$.

If $\mathcal{S} \subset \mathcal{H}$ is a linear subspace and $h \in \mathcal{H}$, let $\Pi(h|\mathcal{S})$ be the projection of h on $\overline{\mathcal{S}}$. For subspaces $\mathcal{S} \subset \bigoplus_{k=1}^d L^2(P_0)$, $\Pi(s|\mathcal{S})$ denotes the (componentwise) projection of a vector $s \in \bigoplus_{k=1}^d L^2(P_0)$ on $\overline{\mathcal{S}}$. Theorem A.4.2 (A) from BKRW directly yields the following result.

Lemma 1 *Assume that the conditions of Theorem 1 hold true. When $m \rightarrow \infty$,*

$$\overline{S}_{\theta_0}^{(m)} = a_1^{(m)}(X^{(1)})g_1 + a_2^{(m)}(X^{(2)})g_2 \longrightarrow \overline{S}_{\theta_0} = \Pi\left(S_{\theta_0}|\mathcal{T}^\perp\right) = \Pi\left(S_{\theta_0}|\overline{H_1 + H_2}\right)$$

in $\bigoplus_{k=1}^d L^2(P_0)$, where $g_i = g(Z, \theta_0)$, $i = 1, 2$.

Let us point out that even if Lemma 1 guarantees the convergence of the iterations $\overline{S}_{\theta_0}^{(m)}$, it is not necessarily true that the sequences $a_1^{(m)}(X^{(1)})g_1$ and $a_2^{(m)}(X^{(2)})g_2$ converge. Sufficient mild conditions are provided in Theorem A.4.2 (C) of BKRW, that are

$$\overline{S}_{\theta_0} = \Pi\left(S_{\theta_0}|\mathcal{T}^\perp\right) = a_1^*(X^{(1)}) \cdot g_1 + a_2^*(X^{(2)}) \cdot g_2 \in \mathcal{T}_1^\perp + \mathcal{T}_2^\perp \quad (6)$$

with $a_1^*(X^{(1)}) \cdot g_1 \in \mathcal{T}_1^\perp \cap (\mathcal{T}_1^\perp \cap \mathcal{T}_2^\perp)^\perp \subset \mathcal{T}_1^\perp$. Moreover, by Proposition A.4.1 of BKRW, condition (6) is equivalent with the existence of a solution $a_1^*g_1$ and $a_2^*g_2$ for the system

$$\begin{cases} a_1^*(X^{(1)}) g_1 = \rho_1 - E\left[a_2^*(X^{(2)}) g_2 g'_1 \mid X^{(1)}\right] V^-(g_1|X^{(1)}) g_1 \\ a_2^*(X^{(2)}) g_2 = \rho_2 - E\left[a_1^*(X^{(1)}) g_1 g'_2 \mid X^{(2)}\right] V^-(g_2|X^{(2)}) g_2, \end{cases} \quad (7)$$

where

$$\begin{aligned} \rho_i = \rho_i(Z, \theta_0) &:= \Pi\left(S_{\theta_0}|\mathcal{T}_i^\perp\right) = E\left(S_{\theta_0}g'_i \mid X^{(i)}\right) V^-(g_i \mid X^{(i)}) g_i \\ &= -\partial_\theta E\left(g'_i \mid X^{(i)}\right) V^-(g_i \mid X^{(i)}) g_i. \end{aligned}$$

(A careful inspection of the proof of Proposition A.4.1 of BKRW shows that the condition $H_1 + H_2 = \mathcal{T}_1^\perp + \mathcal{T}_2^\perp$ is a closed subspace is not necessary for deriving that result, since what is really used in their proof is the relation $H_1^\perp \cap H_2^\perp = (H_1 + H_2)^\perp$. If in addition the system (7) has a unique solution, the backfitting algorithm above is nothing but a convergent iterative procedure for finding it.

In applications, a convenient way to check uniqueness is to prove a contraction property. This is the case for instance if $\mathcal{T}_1^\perp \cap \mathcal{T}_2^\perp = \{0\}$, which in our framework holds if

$$E\left(g_1 g'_2 \mid X^{(1)}, X^{(2)}\right) = 0$$

(in the sequential case, this can be achieved by writing the initial system in an equivalent form satisfying the orthogonal condition above; see subsection 4.1).

In the general case where $\mathcal{T}_1^\perp \cap \mathcal{T}_2^\perp \neq \{0\}$ the system (7) rewritten as in Proposition A.4.1 of BKRW under the form

$$\begin{cases} h_1^* = \Pi(S_{\theta_0} - h_2^* | \mathcal{T}_1^\perp) \\ h_2^* = \Pi(S_{\theta_0} - h_1^* | \mathcal{T}_2^\perp), \end{cases}$$

does not necessarily have the contraction property. In our problem $h_1^* = a_1^* g_1$ and $h_2^* = a_2^* g_2$ with g_1 and g_2 given. Hence it suffices to check a contraction property for $a_1^* g_1$ and $a_2^* g_2$ or some given transformations of them. We will see in subsection 4.2 that in the regression-like models with missing data framework, see Robins, Rotnitzky, Zhao (1994), the equations (7) lead to a contraction property for some given transformations of $a_1^* g_1$ and $a_2^* g_2$.

The ‘backfitting’ algorithm we proposed above involves θ_0 that is unknown. In practice one can use the following steps: (i) build $\tilde{\theta}_n$ a \sqrt{n} -consistent estimator of θ_0 , for instance the *smooth minimum distance estimator (SMD)* like in Lavergne and Patilea (2013); (ii) estimate nonparametrically $a_1^{(m^*)}$ and $a_2^{(m^*)}$ the solution of the ‘backfitting’ algorithm obtained after, say, m^* iterations using $\tilde{\theta}_n$ instead of θ_0 ; and (iii) use classical GMM to construct an efficient estimator $\hat{\theta}^{(m^*)}$ based on the approximate efficient score equations $E(\widehat{S}_\theta) = 0$, where

$$\widehat{S}_\theta = \widehat{a}_1^{(m^*)}(X^{(1)}, \tilde{\theta}_n) g_1(Z, \theta) + \widehat{a}_2^{(m^*)}(X^{(2)}, \tilde{\theta}_n) g_2(Z, \theta),$$

and $\widehat{a}_i^{(m^*)}(X^{(i)}, \tilde{\theta}_n)$ are nonparametric estimates of $a_i^{(m^*)}(X^{(i)}, \theta_0)$, $i = 1, 2$.

4 Applications

In this section we illustrate the utility of our theoretical results for four classes of models: sequential (nested) conditional models, regression-like models with missing data, simultaneous equation models and bivariate binary choice models. The general results in sections 2 and 3 above allow us: (a) to complete a semiparametric efficiency bound result of Chamberlain (1992b); (b) to generalize the mean regression with missing data setting of Robins, Rotnitzky and Zhao (1994) and Tan (2011) to more general moment conditions, which includes for example quantile regressions; (c) to consider a simultaneous equation model with perhaps more appealing orthogonality conditions on the error terms; and (d) to generalize the bivariate probit models.

4.1 Sequential conditional moments

Important cases where equations (7) have an explicit solution are the cases where $\sigma(X^{(1)}) \subset \sigma(X^{(2)})$ holds true. In the case $J = 2$, the model $E(g_j(Z, \theta) | X^{(j)}) = 0$, $j = 1, 2$, defined in (1) can be equivalently written under the form

$$\begin{cases} E(\tilde{g}_1(Z, \theta) | X^{(1)}) = 0 \\ E(g_2(Z, \theta) | X^{(2)}) = 0, \end{cases} \quad (8)$$

where

$$\begin{aligned} \tilde{g}_1(Z, \theta) &= g_1(Z, \theta) \\ &\quad - E(g_1(Z, \theta_0) | X^{(2)}) - V^{-1}(g_2(Z, \theta_0) | X^{(2)}) g_2(Z, \theta). \end{aligned}$$

Here we suppose that $V(g_1(Z, \theta_0) | X^{(1)})$ and $V(g_2(Z, \theta_0) | X^{(2)})$ are invertible and this guarantees that θ_0 is also identified by the equations (8). Recall that g_i is a short notation for $g_i(Z, \theta_0)$ and similarly let \tilde{g}_i replace $\tilde{g}_i(Z, \theta_0)$.

Notice that \tilde{g}_1 is the residual of the projection of g_1 on g_2 with respect to $\sigma(X^{(2)})$ and $E(\tilde{g}_1 | X^{(2)}) = 0$. Let $\tilde{\mathcal{T}}_1$ be the tangent space of the model defined by the first equation in (8). By the definition of \tilde{g}_1 , it is quite clear that condition $\tilde{\mathcal{T}}_1^\perp \cap \mathcal{T}_2^\perp = \{0\}$ holds true. Next, multiplying the i th equation in (7) by g_i , taking conditional expectation given $X^{(i)}$ and finally multiplying by $V^{-1}(g_i | X^{(i)})$, $i = 1, 2$, the system (7) corresponding to model (8) becomes

$$\begin{cases} \tilde{a}_1^*(X^{(1)}) = -\partial_\theta E(\tilde{g}_1' | X^{(1)}) V^{-1}(\tilde{g}_1 | X^{(1)}) \\ \quad - E(\tilde{a}_2^*(X^{(2)}) \cdot g_2 \tilde{g}_1' | X^{(1)}) V^{-1}(\tilde{g}_1 | X^{(1)}) \\ \tilde{a}_2^*(X^{(2)}) = -\partial_\theta E(g_2' | X^{(2)}) V^{-1}(g_2 | X^{(2)}) \\ \quad - E(\tilde{a}_1^*(X^{(1)}) \cdot \tilde{g}_1 g_2' | X^{(2)}) V^{-1}(g_2 | X^{(2)}). \end{cases} \quad (9)$$

Since by definition $E(\tilde{a}_2^*(X^{(2)}) g_2 \tilde{g}_1' | X^{(1)}) = E[\tilde{a}_2^*(X^{(2)}) E(g_2 \tilde{g}_1' | X^{(2)}) | X^{(1)}] = 0$ and $E(\tilde{a}_1^*(X^{(1)}) \tilde{g}_1 g_2' | X^{(2)}) = \tilde{a}_1^*(X^{(1)}) E(\tilde{g}_1 g_2' | X^{(2)}) = 0$ we obtain

$$\begin{cases} \tilde{a}_1^*(X^{(1)}) = -\partial_\theta E(\tilde{g}_1' | X^{(1)}) V^{-1}(\tilde{g}_1 | X^{(1)}) \\ \tilde{a}_2^*(X^{(2)}) = -\partial_\theta E(g_2' | X^{(2)}) V^{-1}(g_2 | X^{(2)}). \end{cases} \quad (10)$$

($\partial_\theta E(\tilde{g}_i')$ denotes the transposed of the matrix $\partial_{\theta'} E(\tilde{g}_i)$.) The efficient score \bar{S}_{θ_0} can then be written as

$$\begin{aligned} \bar{S}_{\theta_0} &= \tilde{a}_1^*(X) \cdot \tilde{g}_1 + \tilde{a}_2^*(X) \cdot g_2 \\ &= -\partial_\theta E(\tilde{g}_1' | X^{(1)}) V^{-1}(\tilde{g}_1 | X^{(1)}) \tilde{g}_1 \\ &\quad - \partial_\theta E(g_2' | X^{(2)}) V^{-1}(g_2 | X^{(2)}) g_2. \end{aligned}$$

In the particular case where $X^{(1)} = X^{(2)} = X$,

$$\begin{aligned}
\bar{S}_{\theta_0} &= \tilde{a}_1^*(X) \cdot \tilde{g}_1 + \tilde{a}_2^*(X) \cdot g_2 \\
&= -\partial_\theta E(\tilde{g}'_1 | X) V^{-1}(\tilde{g}_1 | X^{(1)}) \tilde{g}_1 - \partial_\theta E(g'_2 | X) V^{-1}(g_2 | X) g_2 \\
&= \begin{pmatrix} -\partial_\theta E(\tilde{g}'_1 | X) \\ -\partial_\theta E(g'_2 | X) \end{pmatrix}' V^{-1} \left(\begin{pmatrix} \tilde{g}_1 \\ g_2 \end{pmatrix} | X \right) \begin{pmatrix} \tilde{g}_1 \\ g_2 \end{pmatrix} \\
&= -\partial_\theta E(g' C'(X) | X) V^{-1}(C(X) g | X) C(X) g \\
&= -\partial_\theta E(g' | X) V^{-1}(g | X) g,
\end{aligned}$$

where $g' = (g'_1 \ g'_2)$ and

$$C(X) = \begin{pmatrix} I & -E(g_1 \ g'_2 | X) V^{-1}(g_2 | X) \\ 0 & I \end{pmatrix}$$

is a nonsingular random matrix. This expression of the efficient score directly yields the efficiency bound derived in Chamberlain (1987).

Another important particular case of formulae (10) is provided by models defined by sequential conditional moments; see Chamberlain (1992b), Ai and Chen (2012). Taking $X^{(1)} = X_1$ and $X^{(2)} = (X'_1, X'_2)'$, one obtains

$$\begin{aligned}
\bar{S}_{\theta_0} &= \tilde{a}_1^*(X) \cdot \tilde{g}_1 + \tilde{a}_2^*(X) \cdot g_2 \\
&= -\partial_\theta E(\tilde{g}'_1 | X_1) V^{-1}(\tilde{g}_1 | X_1) \tilde{g}_1 - \partial_\theta E(g'_2 | X_1, X_2) V^{-1}(g_2 | X_1, X_2) g_2.
\end{aligned}$$

Let us point that Chamberlain (1992b) only proves this result for discrete distributions and Ai and Chen (2012) obtain the result in a more general framework (allowing for unknown infinite dimensional parameters in the equations defining the model) but under slightly more restrictive assumptions than in our setting.²

4.2 Regression-like models with missing data

Consider now a regression-like model defined by the equations

$$E[\rho(Y, X^*, \alpha) | X^*] = 0, \tag{11}$$

²In Ai and Chen (2012) it is implicitly required that the class \mathcal{G} appearing in their Assumption A in the Mathematical Appendix is the same for each value of their model parameter α . This variation independent parametrization assumption represents an additional restriction that is unnecessary in our approach. See also van der Laan and Robins (2003), page 18, for some lucid comments on the existence of a variation independent parametrization.

where $\rho(\cdot, \cdot, \cdot)$ is some measurable vector-valued function, α is a (finite-dimension) vector of parameters, and the vector $(Y', X^{*'}) = (Y', X', V')$ is not always completely observed. We also assume that a non-missing indicator δ and some other variable V^0 are always observed. In the following examples we consider two random missingness mechanisms considered respectively by Tan (2011) and Robins, Rotnitzky and Zhao (1994).

Example 1 (i) The vector Y is observed iff $\delta = 1$;

(ii) The vector $W = \begin{pmatrix} X^* \\ V^0 \end{pmatrix}$ is always observed and we have

$$P(\delta = 1 | Y, W) = P(\delta = 1 | W) = \pi(W). \quad (12)$$

Example 2 (i) Let $X^* = \begin{pmatrix} X \\ V \end{pmatrix}$ where X is observed iff $\delta = 1$;

(ii) The vector $W = \begin{pmatrix} Y \\ V \\ V^0 \end{pmatrix}$ is always observed and we have

$$P(\delta = 1 | X, W) = P(\delta = 1 | W) = \pi(W). \quad (13)$$

Let α_0 be the true value of the parameter identified by the model (11). The equation (11) and each of (12) or (13) imply

$$E \left[\frac{\delta}{\pi(W)} \rho(Y, X^*, \alpha_0) | X^* \right] = 0. \quad (14)$$

We can consider this equation at the observational level even for missing X^* , since for missing values of X^* we have $\delta = 0$ which renders the equation noninformative. Note also that (12) and (13) can be written under the unified form

$$P(\delta = 1 | Y, X^*, W) = \pi(W).$$

Therefore, at the observational level, with any of the two examples we obtain a model like

$$\begin{cases} E \left[\frac{\delta}{\pi(W)} \rho(Y, X^*, \alpha_0) | X^* \right] = 0 \\ E \left[\frac{\delta}{\pi(W)} - 1 | W \right] = 0. \end{cases} \quad (15)$$

Moreover, like in Graham (2011), (footnote 8, page 442), it can be shown that, at the observational level, a model given by equation (11) and any of the missing data

mechanism described in Example 1 or Example 2 is equivalent to the model defined by (15).

With our notation, Z is the vector built as the union of all the variables contained in Y, X^*, W and $\delta, \theta = \alpha, g_1(Z, \theta) = \{\delta/\pi(W)\}\rho(Y, X^*, \alpha), g_2(Z, \theta) = \{\delta/\pi(W)\} - 1, X^{(1)} = X^*$ and $X^{(2)} = W$. Let ρ be a short for $\rho(Y, X^*, \alpha_0)$. Then the functions a_1^* and a_2^* defining the efficient score are given by the following equations obtained (see also equations (9)) from equations (7) :

$$\begin{aligned} a_1^*(X^*) &= a_1^*(X^{(1)}) \\ &= -\partial_\alpha E(\rho' | X^*) E^{-1} \left(\frac{1}{\pi(W)} \rho \rho' | X^* \right) \\ &\quad + E \left\{ E[a_1^*(X^*) \rho | W] \frac{1 - \pi(W)}{\pi(W)} \rho' | X^* \right\} E^{-1} \left(\frac{1}{\pi(W)} \rho \rho' | X^* \right); \\ a_2^*(W) &= a_2^*(X^{(2)}) \\ &= E \left[a_1^*(X^*) \rho \frac{\delta}{\pi(W)} \left(\frac{\delta}{\pi(W)} - 1 \right) | W \right] E^{-1} \left[\left(\frac{\delta}{\pi(W)} - 1 \right)^2 | W \right] \\ &= -E[a_1^*(X^*) \rho | W]. \end{aligned}$$

In the particular case where $\rho = \rho(Y, X^*, \alpha_0) = Y - g(X^*, \alpha_0)$ and the selection probability $\pi(W)$ is known, these are exactly the equations obtained in Robins, Rotnitzky and Zhao (1994). They showed that for the regression case, the equation for a_1^* corresponds to a contraction (see the proof of their Proposition 4.2). In subsection 5.3 in the Appendix we show that such a contraction property holds for a more general ρ . Hence we could include in our framework further interesting examples, *e.g.* quantile regressions. The contraction property allows to solve the equations in $a_1^*(X^*)$ and $a_2^*(W)$ by successive approximations.

Let us consider the extended framework where the selection probability is known up to an unknown finite dimension parameter γ_0 , that is

$$P(\delta = 1 | W) = \pi(W, \gamma_0),$$

(see also Robins, Rotnitzky and Zhao (1994), equation (18)). In subsection 5.5 in the Appendix we show that the efficiency score for α_0 has the same expression regardless the selection probability function π is given or depends on the unknown parameter γ_0 . Thus, we extend a result of Robins, Rotnitzky and Zhao (1994), see also Tan (2011), obtained in the particular case of mean regressions.

We close this example with the following remark. Robins, Rotnitzky and Zhao (1994) considered the case where missingness arises only in covariables X^* (that is

also the case considered in our Example 2) and derived the efficient score equations. Tan (2011) obtained formally the same equations with missing regressors *or* missing responses (the case corresponding to our Example 1) using the corresponding definition of W . However, our approach based on conditional moments allows to deeper understand an important difference between the Examples 1 and 2. That is, in the possibly missing responses case we have $\sigma(X^*) \subset \sigma(W)$, so that Example 1 falls in the sequential conditional moments framework where the solutions for a_1^* and a_2^* are explicit. On the other hand, such explicit solutions are *no longer* available in the framework considered by Robins, Rotnitzky and Zhao (1994) and in our Example 2.

4.3 Simultaneous equations models

Consider the linear simultaneous equations

$$\begin{cases} Q = aP + b'X + c'T + \varepsilon \\ P = \alpha Q + \beta'X + \gamma'U + \eta, \end{cases} \quad (16)$$

that could represent, for instance, the equations of a demand (Q) and supply (P) model. Usually it is admitted that the errors ε and η are correlated with the endogenous variables Q and P but are uncorrelated with the exogenous regressors, that means

$$E(\varepsilon | X, T, U) = E(\eta | X, T, U) = 0. \quad (17)$$

Under this strong assumption, the model can be written as

$$E[\rho_1(Z, \theta) | X^{(1)}] = 0, \quad (18)$$

with a two components function ρ_1 , the random vectors $Z = (Q, P, X, T, U)'$, $X^{(1)} = (X, T, U)'$ and parameter $\theta = (a, b', c', \alpha, \beta', \gamma)'$. The model (18) allows for standard inference methods (GMM, instrumental variables, 2SLS) for efficient estimation of θ . With our approach, we are able to relax these assumptions by letting the errors in each equation to be uncorrelated with the corresponding exogenous variables of the same equation, that is we only impose

$$E(\varepsilon | X, T) = 0 \quad \text{and} \quad E(\eta | X, U) = 0, \quad (19)$$

that represent more appealing conditions in the econometric literature. See for instance the equations (14.33) to (14.36) in Wooldridge (2010).

We then obtain the following system of equations of the form given in (1) :

$$\begin{cases} E[\rho_1(Z, \theta) | X^{(1)}] = 0 \\ E[\rho_2(Z, \theta) | X^{(2)}] = 0, \end{cases} \quad (20)$$

where Z and θ are the same as above, but $X^{(1)} = (X, T)'$, $X^{(2)} = (X, U)'$, $\rho_1(Z, \theta) = Q - aP - b'X - c'T$ and $\rho_2(Z, \theta) = P - \alpha Q - \beta'X - \gamma'U$.

Let us point out that, on one hand, one could also take into account additional exogenous instruments for each of the two equations in (16). Let I_1, I_2 denote such instruments. Then one could add the moment equations $E[\rho_1(Z, \theta) I_1] = 0$, $E[\rho_2(Z, \theta) I_2] = 0$ to the system (20). On the other hand, one could also take into account additional observed variables and moment conditions not containing unknown parameters like in Qian and Schmidt (1999), section 4. However, in order to provide a readable illustration of our methodology, for the remainder of this example we stay with the system (16).

If $1 - a\alpha \neq 0$, the system (16) can be written under the restricted reduced form

$$\begin{cases} Q = \delta(b' + a\beta')X + \delta c'T + \delta a\gamma'U + \tilde{\varepsilon} \\ P = \delta(\alpha b' + \beta')X + \delta \alpha c'T + \delta \gamma'U + \tilde{\eta} \end{cases}, \quad (21)$$

where $\delta = (1 - a\alpha)^{-1}$, $\tilde{\varepsilon} = \delta(\varepsilon + a\eta)$ and $\tilde{\eta} = \delta(\alpha\varepsilon + \eta)$. Under the conditions (17), we also have $E(\tilde{\varepsilon} | X, T, U) = E(\tilde{\eta} | X, T, U) = 0$. If the matrix $E(WW')$ has full rank, where $W' = (X', T', U')$, then the coefficients $s_1 = \delta(b + a\beta)$, \dots , $s_6 = \delta\gamma$ are identified and can be consistently estimated by OLS and efficiently estimated in a second stage, using nonparametric estimates of the conditional variances $V(\tilde{\varepsilon} | X, T, U)$ and $V(\tilde{\eta} | X, T, U)$. This two-stage procedure leads to efficient estimators for the vector coefficient θ in the original system (16).³

Suppose now that we only impose the conditions from equations (19). Then $E(\tilde{\varepsilon} | X) = E(\tilde{\eta} | X) = 0$ and one could consistently estimate θ , but the two-stage estimation procedure does no longer yield an efficient estimator. An approximately efficient estimator can be obtained by the approaches proposed in this paper, that is either by increasing the number of instruments for each equation (depending on X and T for the first equation in the system (20) and on X and U for the second one), or by the iterative procedure described in the section 3 above. To illustrate the iterative approach, let us assume the conditional variances $\sigma_{\tilde{\varepsilon}}^2(X, T) = V(\tilde{\varepsilon} | X, T)$ and $\sigma_{\tilde{\eta}}^2(X, U) = V(\tilde{\eta} | X, U)$ are strictly positive. Noting that the functions a_1^* and a_2^* appearing in equations (7) can be arbitrarily defined on the sets $g_1 = 0$ and

³One could also consider the weaker conditions

$$E(\varepsilon | X) = E(\eta | X) = 0, \quad (22)$$

in which case $E(\tilde{\varepsilon} | X) = E(\tilde{\eta} | X) = 0$. A two-stage efficient procedure based on the nonparametric estimation of the conditional variances with respect to X is still possible. However, such weaker conditions like on the errors are less intuitive and justified for modeling purposes. The point we make with this example is that assumptions (19) may appear as the most natural compromise between the 'extreme' conditions (22), the less restrictive, and (19), the most restrictive.

$g_2 = 0$, respectively, we obtain the following system for a_1^* and a_2^* :

$$\begin{cases} a_1^*(X, T) \sigma_\varepsilon^2(X, T) = (E(P | X, T), X', T', 0, 0, 0)' - E[a_2^*(X, U) g_2 g_1 | X, T] \\ a_2^*(X, U) \sigma_\eta^2(X, U) = (0, 0, 0, E(Q | X, U), X', U')' - E[a_1^*(X, T) g_1 g_2 | X, U] \end{cases}.$$

4.4 Generalization of bivariate probit models

Consider the following system involving two binary outcomes Y_1 and Y_2 :

$$\begin{cases} Y_1 = \mathbf{1} \{X_1^T \beta_1 - \varepsilon_1 > 0\} \\ Y_2 = \mathbf{1} \{X_2^T \beta_2 - \varepsilon_2 > 0\}. \end{cases} \quad (23)$$

If $\varepsilon = (\varepsilon_1, \varepsilon_2)$ is bivariate normal and independent of (X_1, X_2) this is the classical bivariate probit model (see, for example, Wooldridge (2010), page 595).

Assuming only that ε_i is independent of X_i and the distribution function F_{ε_i} of ε_i is known, for $i \in \{1, 2\}$, then the model becomes

$$\begin{cases} E[Y_1 - F_{\varepsilon_1}(X_1^T \beta_1) | X_1] = 0 \\ E[Y_2 - F_{\varepsilon_2}(X_2^T \beta_2) | X_2] = 0. \end{cases} \quad (24)$$

Since the joint law of $(\varepsilon_1, \varepsilon_2)$ is not necessarily known, as in the classical bivariate probit model, a maximum likelihood approach to estimate the parameters in (23) is no longer available, but one can (approximately) efficiently estimate the parameters with our approach. Note that in this form (24) no specific assumption is needed on the joint law of $(\varepsilon_1, \varepsilon_2)$. In particular we do not need to consider the nuisance parameter $\rho = cov(\varepsilon_1, \varepsilon_2)$ that is usually integrated in the gaussian likelihood in order to estimate β_1 and β_2 . This approach still works if ε_i and X_i are not independent, assuming instead that the conditional distribution function $F_{\varepsilon_i}(\cdot, x_i)$ of ε_i given that $X_i = x_i$ is known, for $i \in \{1, 2\}$.

References

- [1] Ahn, S. C., Schmidt, P. (1999): Estimation of linear panel data models using GMM. In: Mátyás, L. (ed.) *Generalized Method of Moments Estimation*. Cambridge University Press
- [2] Ai, C., Chen, X. (2003): Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica* 71, 1795-1843

- [3] Ai, C., Chen, X. (2012): The semiparametric efficiency bound for models of sequential moment restrictions containing unknown functions. *Journal of Econometrics* 170, 442-457
- [4] Bickel, P. J., Klaassen, C. A. J., Ritov, Y., Wellner, J. A. (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. The John Hopkins University Press
- [5] Chamberlain, G. (1987): Asymptotic efficiency in estimation with conditional moment restrictions. *Econometrica* 34, 305-334
- [6] Chamberlain, G. (1992a): Efficiency bounds for semiparametric regression. *Econometrica* 60, 567-596
- [7] Chamberlain, G. (1992b): Comment: Sequential moment restrictions in panel data. *Journal of Business & Economic Statistics* 10, 20-26
- [8] Chen, X., Pouzo, D. (2009): Efficient estimation of semiparametric conditional moment models with possibly nonsmooth residuals. *Journal of Econometrics* 152, 46-60
- [9] Gallant, A. R. (1975): Seemingly unrelated nonlinear regressions. *Journal of Econometrics* 3, 35-50)
- [10] Graham, B. (2011): Efficiency bounds for missing data models with semiparametric restrictions. *Econometrica* 79, 437-452
- [11] Hahn, J. (1997): Efficient estimation of panel data models with sequential moment restrictions. *Journal of Econometrics* 79, 1-21
- [12] Han, C., Phillips, P. C. B. (2006): GMM with many moment conditions. *Econometrica* 74, 147-192
- [13] Hansen, L. P. (1982): Large sample properties of generalized method of moments estimators. *Econometrica* 50, 1029-1054
- [14] Hansen, L. P. (2008): Generalized Method of Moments. In : Durlauf, S. N., Blume, L. E. (eds.) *The New Palgrave Dictionary of Economics*. Second Edition. Palgrave Macmillan
- [15] Hansen, L. P., Sargent, T. J. (1991): *Rational Expectations Econometrics*. Westview Press
- [16] Ibragimov, I. A., Has'minskii, R. Z. (1981): *Statistical Estimation. Asymptotic Theory*. Springer-Verlag, New-York

- [17] Jun, S., Pinske, J. (2009): Efficient semiparametric seemingly unrelated quantile regression estimation. *Econometric Theory* 25, 1392-1414
- [18] Lavergne, P. (2008): A Cauchy-Schwarz inequality for expectation of matrices. <http://econpapers.repec.org/RePEc:sfu:sfudps:dp08-07>
- [19] Lavergne, P., Patilea, V. (2013): Smooth minimum distance estimation and testing with conditional estimating equations: uniform in bandwidth theory. *Journal of Econometrics*, 177, 47-59
- [20] Müller, U. (2009): Estimating linear functionals in nonlinear regression with responses missing at random. *Annals of Statistics* 37, 2245-2277
- [21] Müller, U.U., Van Keilegom, I. (2012): Efficient parameter estimation in regression with missing responses. *Electronic Journal of Statistics* 6, 1200-1219
- [22] Newey, W. K. (1990): Semiparametric efficiency bounds. *Journal of Applied Econometrics* 5, 99-135
- [23] Newey, W. K. (1993): Efficient estimation of models with conditional moment restrictions. In : *Handbook of statistics*, 11, 419-454.
- [24] Newey, W. K. (2001): Conditional moment restrictions in censored and truncated regression models. *Econometric Theory* 17, 863-888
- [25] Newey, W. K. (2004): Efficient semiparametric estimation via moment restrictions. *Econometrica* 72, 1877-1897
- [26] Qian, H., Schmidt, P. (1999): Improved instrumental variables and generalized method of moments estimators. *Journal of Econometrics*, 91, 145-169
- [27] Robins, J. M., Rotnitzky A., Zhao L. P. (1994): Estimation of regression coefficients when some regressors are not always observed. *Journal of the American Statistical Association* 89, 846-866
- [28] Tan, Z. (2011): Efficient restricted estimators for conditional mean models with missing data. *Biometrika* 98, 663-684
- [29] Tsiatis, A. A. (2006): *Semiparametric Theory and Missing Data*. Springer, New York
- [30] van der Laan, M. J., Robins, J. M. (2003): *Unified Methods for Censored Longitudinal Data and Causality*. Springer-Verlag, New York
- [31] van der Vaart, A.W. (1998): *Asymptotic Statistics*. Cambridge University Press

- [32] Wei, Y., Ma, Y., Carroll, R. J. (2012): Multiple imputation in quantile regression. *Biometrika* 99, 423-438
- [33] Wooldridge, J.M. (2010): *Econometric Analysis of Cross Section and Panel Data* (2e). MIT Press

5 Appendix

For a model \mathcal{P} (resp. \mathcal{P}_j) and a probability measure P in the model, let $\dot{\mathcal{P}}_P$ (resp. $\dot{\mathcal{P}}_{j,P}$) denote the tangent cone of the model \mathcal{P} (resp. \mathcal{P}_j) at P . When there is no possible confusion, we simply write $\dot{\mathcal{P}}_P$ (resp. $\dot{\mathcal{P}}_{j,P}$). Let $\mathcal{T}(\mathcal{P}, P)$ denote the tangent space of a model \mathcal{P} at some probability measure $P \in \mathcal{P}$, that means the closure of the linear span of the tangent set $\dot{\mathcal{P}}_P$. By definition, both the tangent cone and the tangent space are subsets of $L^2(P)$.

5.1 A general lemma

The following result is a generalization of Theorem 1 in Newey (2004) where only the case of models defined by conditional moments, with the same conditioning variables, was considered.

Lemma 2 *Let $P_0 \in \mathcal{P} \subset \mathcal{P}_1$ be the true law of the vector $Z \in \mathcal{Z}$ and $\theta_0 = \psi(P_0)$ for a map $\psi : \mathcal{P}_1 \rightarrow \mathbb{R}^d$ differentiable at P_0 relative to the tangent cone $\dot{\mathcal{P}}_{1,P_0}$. Let $\{\mathcal{P}_k\}_{k \in \mathbb{N}^*}$ be a decreasing family of statistical models such that*

$$\mathcal{P}_1 \supset \mathcal{P}_2 \supset \dots \supset \mathcal{P}_j \supset \mathcal{P}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{P}_k \supset \mathcal{P} \ni P_0 \quad (25)$$

and

$$\bigcap_{k=1}^{\infty} \mathcal{T}_k = \mathcal{T}, \quad (26)$$

where $\mathcal{T} = \mathcal{T}(\mathcal{P}, P_0)$ and $\mathcal{T}_k = \mathcal{T}(\mathcal{P}_k, P_0)$, $k \in \mathbb{N}^*$. Then

$$I_{\theta_0}(\mathcal{P}) = \lim_{k \rightarrow \infty} I_{\theta_0}(\mathcal{P}_k),$$

where $I_{\theta_0}(\mathcal{P})$ stands for the Fisher information on $\theta_0 = \psi(P_0)$ in the model \mathcal{P} .

For the definition of the Fisher information $I_{\theta_0}(\mathcal{P})$ on $\theta_0 = \psi(P_0)$ in the model \mathcal{P} we refer to Bickel, Klaassen, Ritov and Wellner (1993) or van der Vaart (1998); see also Newey (1990). When the models \mathcal{P}_k , $k \in \mathbb{N}^*$, are defined by an increasing number of moment conditions with the same conditioning vectors, condition (26) is exactly the so-called spanning condition of Newey (2004).

Proof of Lemma 2. By definition (van der Vaart (1998), p. 363), there exists a continuous linear map $\dot{\psi} : L^2(P_0) \rightarrow \mathbb{R}^d$ such that for any $g \in \dot{\mathcal{P}}_{P_0} \subset L^2(P_0)$ and a submodel $(-\varepsilon, \varepsilon) \ni t \mapsto P_t$ with score function g ,

$$\frac{\psi(P_t) - \psi(P_0)}{t} \xrightarrow[t \rightarrow 0]{} \dot{\psi}(g).$$

By the Riesz representation theorem, there exists a unique d -dimension vector-valued function having the components in $L^2(P_0)$ such that $\dot{\psi}(h) = E_{P_0}(\bar{\psi}h)$ for every $h \in L^2(P_0)$. In particular,

$$\dot{\psi}(g) = E_{P_0}(\bar{\psi}g) = \int \bar{\psi}gdP_0, \quad \forall g \in \dot{\mathcal{P}}_{P_0} \subset L^2(P_0).$$

Let $\tilde{\psi}$ and $\tilde{\psi}_k$ denote the elements of $[L^2(P_0)]^d$ obtained by componentwise projections of $\bar{\psi}$ on the tangent spaces $\mathcal{T} \subset L^2(P_0)$ and $\mathcal{T}_k \subset L^2(P_0)$, respectively. The Fisher information matrices on $\theta_0 = \psi(P_0)$ in the models \mathcal{P} , \mathcal{P}_k at P_0 are then defined by

$$I_{\theta_0}^{-1}(\mathcal{P}) = V_{P_0}(\tilde{\psi}) = E_{P_0}(\tilde{\psi}\tilde{\psi}'), \quad I_{\theta_0}^{-1}(\mathcal{P}_k) = V_{P_0}(\tilde{\psi}_k), \quad k \in \mathbb{N}^*.$$

From

$$\mathcal{P}_1 \supset \mathcal{P}_2 \supset \dots \supset \mathcal{P}_k \supset \mathcal{P}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{P}_k \supset \mathcal{P}$$

we deduce that

$$\dot{\mathcal{P}}_1 \supset \dot{\mathcal{P}}_2 \supset \dots \supset \dot{\mathcal{P}}_k \supset \dot{\mathcal{P}}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \dot{\mathcal{P}}_k \supset \dot{\mathcal{P}},$$

and

$$\mathcal{T}_1 \supset \mathcal{T}_2 \supset \dots \supset \mathcal{T}_k \supset \mathcal{T}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{T}_k = \mathcal{T},$$

where the last equality is due to (4). By Lemma 4.5 of Hansen and Sargent (1991),

$$\lim_{k \rightarrow \infty} I_{\theta_0}^{-1}(\mathcal{P}_k) = \lim_{k \rightarrow \infty} V_{P_0}(\Pi(\bar{\psi}|\mathcal{T}_k)) = V_{P_0}(\Pi(\bar{\psi}|\mathcal{T})) = V_{P_0}(\tilde{\psi}) = I_{\theta_0}^{-1}(\mathcal{P}).$$

■

Remark 1 Even if $\bigcap_{k=1}^{\infty} \mathcal{P}_k = \mathcal{P}$, condition (26) is not necessarily fulfilled. To see this, consider a symmetric density f_0 on the real line and let s_1, s_2 be two odd functions such that $|s_1|, |s_2| \leq 1$ (e.g. $s_l(x) = x^{2l-1} \mathbf{1}\{|x| \leq 1\}$, $l = 1, 2$). For any $k \in \mathbb{N}^*$ and $t \in [-1, 1]$, define

$$f_t(x) = f_0(x)[1 + t s_2(x)], \quad f_{t,k}(x) = k f_0(kx)[1 + t s_1(x)]$$

and consider the following models defined by their densities with respect to $\lambda_{\mathbb{R}}$ the Lebesgue measure on the real line : $\mathcal{Q}_k = \{f_{t,k} \cdot \lambda_{\mathbb{R}} : t \in [-1, 1]\}$, $k \in \mathbb{N}^*$, and

$$\mathcal{P} = \{f_t \cdot \lambda_{\mathbb{R}} : t \in [-1, 1]\}, \quad \mathcal{P}_k = \mathcal{P} \cup \bigcup_{m=k}^{\infty} \mathcal{Q}_m, \quad k \in \mathbb{N}^*.$$

Then we have

$$\mathcal{P}_1 \supset \mathcal{P}_2 \supset \dots \supset \mathcal{P}_k \supset \mathcal{P}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{P}_k = \mathcal{P}.$$

To describe the corresponding tangent spaces, notice that

$$\forall k \geq 1, \quad \partial_t \log f_{t,k}(x)|_{t=0} = s_1(x) \quad \text{and} \quad \partial_t \log f_t(x)|_{t=0} = s_2(x),$$

and thus $\dot{\mathcal{P}} = \{a s_2(x) : a \in \mathbb{R}\}$,

$$\dot{\mathcal{P}}_k = \{a s_2(x) : a \in \mathbb{R}\} \cup \{b s_1(x) : b \in \mathbb{R}\}, \quad k \in \mathbb{N}^*.$$

Then $\mathcal{T} = \{a s_2(x) : a \in \mathbb{R}\}$,

$$\mathcal{T}_k = \{a s_2(x) + b s_1(x) : a, b \in \mathbb{R}\}, \quad k \in \mathbb{N}^*.$$

This shows that

$$\bigcap_{k=1}^{\infty} \dot{\mathcal{P}}_k \not\supseteq \dot{\mathcal{P}} \quad \text{and} \quad \bigcap_{k=1}^{\infty} \mathcal{T}_k \not\supseteq \mathcal{T},$$

even if the decreasing sequence of models $\{\mathcal{P}_k\}_{k \in \mathbb{N}^*}$ is such that $\bigcap_{k=1}^{\infty} \mathcal{P}_k = \mathcal{P}$.

5.2 Comments on the Assumptions T and SP

Let $\varphi_1(Z), \dots, \varphi_J(Z)$ be bounded vector-valued functions of dimensions d_1, \dots, d_J , respectively. For each $j, k \in \{1, \dots, J\}$ we search a $p_j \times d_k$ -matrice $\alpha_k^{(j)}(\underline{X})$ with bounded entries such that

$$b_j(Z) = \alpha_1^{(j)}(\underline{X}) \varphi_1(Z) + \dots + \alpha_J^{(j)}(\underline{X}) \varphi_J(Z) \in \mathbb{R}^{p_j}, \quad j \in \{1, \dots, J\},$$

satisfy a strong form of Assumption T.1, that is

$$E(b_j g'_i | \underline{X}) = 0, \quad \forall i, j \in \{1, \dots, J\}, i \neq j.$$

Recall that for each $i \in \{1, \dots, J\}$, g_i takes values in \mathbb{R}^{p_i} . Then one has the following system for $\alpha_1^{(j)}, \dots, \alpha_J^{(j)}$:

$$\sum_{k=1}^J \alpha_k^{(j)} E(\varphi_k g'_i | \underline{X}) = 0, \quad \forall i, j \in \{1, \dots, J\}, i \neq j. \quad (27)$$

For a given $j \in \{1, \dots, J\}$, there are $p_j \times \sum_{i=1, i \neq j}^J p_j$ equations with $p_j \times \sum_{k=1}^J d_k$ unknown functions (the components of $\alpha_k^{(j)}$'s). If $d_k = p_k$ and the components of $\alpha_k^{(k)}$ are taken constant (for example $\alpha_k^{(k)} \equiv 1$), $1 \leq k \leq J$, then one obtains a unique solution for $\alpha_1^{(j)}, \dots, \alpha_J^{(j)}$ from equations (27), provided a suitable choice of the φ_k 's. The functions g_1, \dots, g_J defining the model would be natural candidates for the φ_k 's if $d_k = p_k$. However, the functions g_1, \dots, g_J are not necessarily bounded and some adjustments are required. For a subset $A \subset \text{supp}Z$, we use the following notations : $g_{i,A} = g_i(Z, \theta_0) \mathbf{1}\{Z \in A\}$, $i = 1, \dots, J$. If we take the subset A of the support of Z such that $g_{1,A}, \dots, g_{J,A}$ are bounded, we can choose $\varphi_k = g_{k,A}$, $k = 1, \dots, J$. Noting that

$$E(g_{j,A} g'_i | \underline{X}) = E(g_{j,A} g'_{i,A} | \underline{X}), \quad \forall i, j \in \{1, \dots, J\},$$

with a suitable choice of A , one can take

$$b_j = \alpha_1^{(j)} g_{1,A} + \dots + g_{j,A} + \dots + \alpha_J^{(j)} g_{J,A} \in \mathbb{R}^{p_j}, \quad j \in \{1, \dots, J\}, \quad (28)$$

in Assumption T. Here $\alpha_k^{(j)} = \alpha_k^{(j)}(\underline{X})$ are the solutions of the system

$$\sum_{k=1}^J \alpha_k^{(j)} E(g_{k,A} g'_{i,A} | \underline{X}) = 0, \quad \forall i, j \in \{1, \dots, J\}, i \neq j, \quad (29)$$

where $\alpha_j^{(j)} = (1, \dots, 1)' \in \mathbb{R}^{p_j}$ for $j \in \{1, \dots, J\}$. With this particular choice of functions b_j , Assumption T can be replaced by the following assumption.

Assumption T' There exist a subset $A \subset \text{supp}Z$ such that the functions b_1, \dots, b_J defined by (28) and (29) satisfy Assumptions T.2 and T.3.

Let us detail the construction of the functions b_1, \dots, b_J in the particular case $J = 2$. Define

$$b_i = g_{i,A} - E(g_{i,A} g'_{j,A} | X^{(1)}, X^{(2)}) E^{-1}(g_{j,A} g'_{j,A} | X^{(1)}, X^{(2)}) g_{j,A}, \quad (30)$$

for $(i, j) \in \{(1, 2), (2, 1)\}$ where $E^{-1}(g_{j,A} g'_{j,A} | X^{(1)}, X^{(2)})$ stands for the inverse of the matrix $E(g_{j,A} g'_{j,A} | X^{(1)}, X^{(2)})$ that is supposed to exist. Assumption T could then be replaced by the following assumption.

Assumption T'' (The case $J = 2$). There exist a subset $A \subset \text{supp}Z$ such that the functions $g_{1,A}, g_{2,A}$ are bounded and b_1, b_2 defined by (30) satisfy Assumptions T.2 and T.3 and, for $i \in \{1, 2\}$, $E(g_{i,A} g'_{i,A} | X^{(1)}, X^{(2)})$ is invertible (a.s.) and

$$\left\| E^{-1}(g_{i,A} g'_{i,A} | X^{(1)}, X^{(2)}) \right\|_{\infty} < \infty \quad (a.s.).$$

If the functions g_1 and g_2 are bounded, Assumption T''.1 with $A = \text{supp}Z$ is slightly stronger than usual assumptions appearing in the literature (see, for example, Assumptions 1 and 2 of Newey (2004), Assumptions 2 and 3 of Ai and Chen (2012)). Meanwhile, our results are derived under Assumption T which is also more general than Assumption T'', is closely related to the usual regularity conditions appearing in the literature, but not exactly comparable with them. Let us also point out that Assumption T.1 is not required in the case $J = 1$ where one can simply take $b_1 = g_{1,A}$.

Note that Assumption SP.1 does not necessarily mean that the parameters θ and η are completely separated, as it is sometimes implicitly assumed in the literature. In fact θ and η are connected since the functional parameter η can have θ among its arguments. Assumption SP only means that when considering the density of $P_{\theta, \eta}$ with respect to a dominating measure μ we could write it under the form

$$f(\cdot, \theta, \eta(v(\cdot, \theta))),$$

with f and v having a known form, where $f(\cdot, \theta_0, \eta(v(\cdot, \theta_0)))$ and $f(\cdot, \theta, \eta_0(v(\cdot, \theta)))$ belong to the model \mathcal{P} for every $\theta \in \Theta$ and $\eta \in H$. For example, in the conditional mean setting with one conditioning vector

$$E[Y - m(X, \theta) | X] = 0,$$

we can take H as the set of zero conditional mean densities of $Z = (Y', X)'$, i.e.

$$H = \left\{ p(y, x) \cdot \gamma(x) : p \geq 0, \gamma \geq 0, \int p(y, x) dy = 1, \int yp(y, x) dy = 0, \forall x, \right. \\ \left. \int \gamma(x) dx = 1 \right\}$$

and $v(y, x, \theta) = (y - m(x, \theta), x)$, so that

$$\eta(v(z, \theta)) = \eta(y - m(x, \theta), x) = p(y - m(x, \theta), x) \cdot \gamma(x)$$

and

$$f(z, \theta, \eta(v(z, \theta))) = \eta(v(z, \theta)).$$

In the proof of Theorem 1 we identify the density $f(\cdot, \theta, \eta(v(\cdot, \theta)))$ with the infinite dimensional nuisance parameter η which is itself a density.

5.3 Proof of Theorem 1

For any $k \in \mathbb{N}^*$, let \mathcal{P}_k be the model defined by equation (4) and \mathcal{P} the model defined by equation (1). Then

$$\mathcal{P}_1 \supset \mathcal{P}_2 \supset \dots \supset \mathcal{P}_k \supset \mathcal{P}_{k+1} \supset \dots \supset \bigcap_{k=1}^{\infty} \mathcal{P}_k = \mathcal{P}.$$

Hence the stated result is a direct consequence of Lemma 2, provided that condition (26) holds for the tangent spaces of \mathcal{P} and \mathcal{P}_k , $k \in \mathbb{N}^*$, at θ_0 .

For each $j \in \{1, \dots, J\}$, any $z \in \mathcal{Z} \subset \mathbb{R}^q$ could be partitioned in two subvectors $y^{(j)} \in \mathbb{R}^{q-q_j}$ and $x^{(j)} \in \mathbb{R}^{q_j}$ with $x^{(j)}$ in the support of $X^{(j)}$. Recall that $P_{X^{(j)}}$ denotes the law of $X^{(j)}$. Model \mathcal{P} is then defined by the set of conditions

$$\int g_j(z, \theta) f(z, \theta) dy^{(j)} = 0 \quad P_{X^{(j)}} - a.s., \quad j \in \{1, \dots, J\}; \quad (31)$$

for a fixed k , the model \mathcal{P}_k is defined by

$$\int g_j(z, \theta) f(z, \theta) w_r^{(j)}(x^{(j)}) dz = 0, \quad j \in \{1, \dots, J\}, \quad r \in \{1, \dots, k\}. \quad (32)$$

Consider now a regular parametric family $\{f_t\}_{t \in (-\varepsilon, \varepsilon)}$ of densities satisfying (31), that means that there exist parameters $\theta_t \in \Theta$ such that, for any $t \in (-\varepsilon, \varepsilon)$ and $P_{X^{(j)}} - a.s.$,

$$\int g_j(z, \theta_t) f_t(z, \theta_t) dy^{(j)} = 0, \quad \forall j \in \{1, \dots, J\}. \quad (33)$$

Let

$$\begin{aligned}\dot{\theta} &= \left. \frac{\partial \theta_t}{\partial t} \right|_{t=0}, \\ s &= \left. \partial_t \log f_t(Z, \theta_t) \right|_{t=0}, \quad S_{\theta_0} = \left. \partial_\theta \log f(Z, \theta) \right|_{\theta=\theta_0}, \\ s_1 &= \left. \partial_t \log f_t(Z, \theta_t) \right|_{t=0} = s + S'_{\theta_0} \dot{\theta}.\end{aligned}$$

Here and in the following, the derivatives of the log-densities are to be understood in the mean square sense, see Ibragimov and Has'minskii (1981), page 64. Differentiating with respect to t in (33) we obtain

$$\partial_{\theta'} E \left[g_j(Z, \theta_0) | X^{(j)} \right] \dot{\theta} + E \left[g_j(Z, \theta_0) s_1(Z) | X^{(j)} \right] = 0, \quad (34)$$

$\forall j \in \{1, \dots, J\}$. Since $\dot{\theta} \in \mathbb{R}^d$ could be arbitrary, we deduce that

$$\begin{aligned}\partial_{\theta'} E \left[g_j(Z, \theta_0) | X^{(j)} \right] + E \left[g_j(Z, \theta_0) S'_{\theta_0}(Z) | X^{(j)} \right] &= 0, \\ E \left[g_j(Z, \theta_0) s(Z) | X^{(j)} \right] &= 0,\end{aligned}$$

for each $j \in \{1, \dots, J\}$. The last equation and the expression of the score functions s_1 suggest a tangent space $\mathcal{T} = \mathcal{T}(P, P_0)$ of the form

$$\mathcal{T} = \overline{\text{lin}} S_{\theta_0} + \left\{ s : E(s^2) < \infty, E(s) = 0, E \left[g_j(Z, \theta_0) s(Z) | X^{(j)} \right] = 0, 1 \leq j \leq J \right\}. \quad (35)$$

On the other hand, the tangent space $\mathcal{T}_k = \mathcal{T}(P_k, P_0)$ corresponding to the model defined by the equations (32) is given by vectors satisfying the unconditional moment equations

$$\partial_{\theta'} E \left[g_j(Z, \theta_0) w_r^{(j)}(X^{(j)}) \right] \dot{\theta} + E \left[g_j(Z, \theta_0) s_1(Z) w_r^{(j)}(X^{(j)}) \right] = 0, \quad (36)$$

$1 \leq j \leq J, 1 \leq r \leq k$. This yields the tangent spaces

$$\begin{aligned}\mathcal{T}_k &= \overline{\text{lin}} S_{\theta_0} + \left\{ s : E(s^2) < \infty, E(s) = 0, E \left[g_j(Z, \theta_0) s(Z) w_r^{(j)}(X^{(j)}) \right] = 0, \right. \\ &\quad \left. \forall 1 \leq j \leq J, \forall 1 \leq r \leq k \right\};\end{aligned}$$

see for instance Example 3, section 3.2 in Bickel, Klaassen, Ritov and Wellner (1993). Since the functions $w_r^{(j)}(X^{(j)})$, $r \in \mathbb{N}^*$, span $L^2(P_{X^{(j)}})$, equations (34) are satisfied

if and only if equations (36) are satisfied for any $k \in \mathbb{N}^*$. In other words, the equivalent of the spanning condition of Newey (2004), see our equation (26) above, is satisfied and we can apply Lemma 2 to conclude that $I_{\theta_0} = \lim_{k \rightarrow \infty} I_{\theta_0}^{(k)}$.

The proof will be complete if we show that the tangent space $\mathcal{T} = \mathcal{T}(\mathcal{P}, P_0)$ is indeed the set described in equation (35). It is quite easy to see that equations (34) guarantee the inclusion “ \subset ” in display (35). To show the reverse inclusion, since $\overline{\text{lin}} S_{\theta_0} \subset \mathcal{T}$ by the definition of the tangent space, it suffices to prove that $\mathcal{T}' \subset \mathcal{T}$, where

$$\begin{aligned} \mathcal{T}' &= \mathcal{T}'(\mathcal{P}, P_0) \\ &= \left\{ s : E(s^2) < \infty, E(s) = 0, E \left[g_j(Z, \theta_0) s(Z) \mid X^{(j)} \right] = 0, 1 \leq j \leq J \right\}. \end{aligned}$$

Let f_0 denote the true density of the vector Z . Take $s \in \mathcal{T}'$ and suppose for the moment that s is bounded. Then, for real numbers t with sufficiently small absolute values, the functions $f_t = (1 + t \cdot s) f_0$ are densities on \mathcal{Z} and if E_{f_t} denotes expectation with respect to the law defined by f_t ,

$$\begin{aligned} E_{f_t} \left[g_j(Z, \theta_0) a \left(X^{(j)} \right) \right] &= E \left[g_j(Z, \theta_0) a \left(X^{(j)} \right) \right] \\ &\quad + t E \left[g_j(Z, \theta_0) s(Z) a \left(X^{(j)} \right) \right] = 0, \end{aligned}$$

for any square-integrable function $a(X^{(j)})$, so that $E_{f_t}[g_j(Z, \theta_0) | X^{(j)}] = 0, 1 \leq j \leq J$. Moreover,

$$\partial_t \log f_t|_{t=0} = \partial_t \log(1 + t \cdot s)|_{t=0} = s,$$

which means that the family of densities $\{f_t\}_{|t| < \varepsilon}$ defines a submodel of model (1) for which the tangent vector at $t = 0$ is exactly s . Next, we have to extend the argument to unbounded functions s . If $\mathcal{M} \subset L^2(P_0)$ is the subspace of bounded functions of Z , it remains to show that $\mathcal{M} \cap \mathcal{T}'$ is dense in \mathcal{T}' . One may consider this step obvious since any unbounded square integrable function can be approximated by a sequence of bounded functions, see for instance Ai and Chen (2003), page 1838. We argue that this well-known approximation result cannot be directly applied to our context, as it is also the case in other contexts considered in the efficiency bounds literature. Indeed, here we are in the following situation: we have J infinite-dimension closed subspaces $\mathcal{T}'_1, \dots, \mathcal{T}'_J$ such that $\mathcal{T}' = \mathcal{T}'_1 \cap \dots \cap \mathcal{T}'_J$, $\overline{\mathcal{M} \cap \mathcal{T}'_j} = \mathcal{T}'_j, 1 \leq j \leq J$, and we need that $\overline{\mathcal{M} \cap \mathcal{T}'} = \mathcal{T}'$. To our best knowledge, there is no general mathematical result which would allow us to claim that $\mathcal{M} \cap \mathcal{T}'$ is dense in \mathcal{T}' without any further argument. That is why we have to provide a proof adapted to the case we consider herein, for which our Assumption T is well suited.

Herein, the norm of a vector (or matrix) should be understood as the sum of componentwise norms. Since \mathcal{M} is dense in $L^2(P_0)$, for a fixed $s \in \mathcal{T}'$ there exist a sequence $\{t_n\}_n \subset \mathcal{M}$ such that

$$\|s - t_n\|_{L^2(P_0)} \xrightarrow{n \rightarrow \infty} 0 \quad \text{and} \quad \max_{1 \leq k \leq J} \left\| E \left(t_n g'_k \mid X^{(k)} \right) \right\|_{\infty} < \infty \quad (a.s.) \quad (37)$$

Define

$$u_n = t_n - \sum_{k=1}^J E \left(t_n g'_k \mid X^{(k)} \right) E^{-1} \left(b_k g'_k \mid X^{(k)} \right) b_k.$$

Since $\{t_n\}_n \subset \mathcal{M}$ satisfies the condition (37), by Assumption T.2 and T.3 we have, for each $k \in \{1, \dots, J\}$,

$$\left\| E \left(t_n g_k \mid X^{(k)} \right) E^{-1} \left(b_k g'_k \mid X^{(k)} \right) b_k \right\|_{\infty} < \infty \quad (a.s.),$$

and thus $u_n \in \mathcal{M}$. Then, for $j \in \{1, \dots, J\}$, using Assumption T.1,

$$\begin{aligned} & E \left(g_j u'_n \mid X^{(j)} \right) \\ &= E \left(g_j t'_n \mid X^{(j)} \right) - \sum_{k=1}^J E \left[g_j b'_k E^{-1} \left(g_k b'_k \mid X^{(k)} \right) E \left(g_k t'_n \mid X^{(k)} \right) \mid X^{(j)} \right] \\ &= E \left(g_j t'_n \mid X^{(j)} \right) \\ &\quad - \sum_{k=1}^J E \left[\underbrace{E \left(g_j b'_k \mid X^{(j)}, X^{(k)} \right)}_{=0 \text{ if } j \neq k} E^{-1} \left(g_k b'_k \mid X^{(k)} \right) E \left(g_k t'_n \mid X^{(k)} \right) \mid X^{(j)} \right] \\ &= \underline{\underline{E \left(g_j t'_n \mid X^{(j)} \right) - E \left(g_j b'_j \mid X^{(j)} \right) E^{-1} \left(g_j b'_j \mid X^{(j)} \right) E \left(g_j t'_n \mid X^{(j)} \right)}} \\ &= \underline{\underline{0}}. \end{aligned} \quad (38)$$

Moreover,

$$\begin{aligned} s - u_n &= s - t_n + t_n - u_n \\ &= s - t_n + \sum_{k=1}^J E \left[(t_n - s) g'_k \mid X^{(k)} \right] E^{-1} \left(b_k g'_k \mid X^{(k)} \right) b_k, \end{aligned}$$

which, by Assumption T.2, entails

$$\begin{aligned} & \|s - u_n\|_{L^2(P_0)} \\ & \leq \|s - t_n\|_{L^2(P_0)} + C \cdot c_b \sum_{k=1}^J \|b_k\|_{\infty} \left\| E \left[(t_n - s) g'_k \mid X^{(k)} \right] \cdot \mathbf{1}\{b_k(Z) \neq 0\} \right\|_{L^2(P_0)}, \end{aligned}$$

for some constant C depending on d_1, \dots, d_J . Let

$$\delta_k^2(X^{(k)}) = E(\mathbf{1}\{b_k(Z) \neq 0\} \mid X^{(k)}).$$

Noting that, for $k \in \{1, \dots, J\}$,

$$\begin{aligned} & \left\| E \left[(t_n - s) g_k \mid X^{(k)} \right] \cdot \mathbf{1}\{b_k \neq 0\} \right\|_{L^2(P_0)}^2 \\ &= E \left\{ E \left[(t_n - s) \delta_k g'_k \mid X^{(k)} \right] \cdot E \left[(t_n - s) \delta_k g_k \mid X^{(k)} \right] \right\} \\ \text{(Cauchy - Schwarz)} &\leq E \left\{ E \left[(t_n - s)^2 \mid X^{(k)} \right] E \left(\delta_k^2 g'_k g_k \mid X^{(k)} \right) \right\} \\ \text{(Assumption T.3)} &\leq \|t_n - s\|_{L^2(P_0)}^2 \left\| E \left(g'_k g_k \mid X^{(k)} \right) \delta_k^2 \right\|_{\infty}^2, \end{aligned}$$

we finally obtain $\|s - u_n\|_{L^2(P_0)} \rightarrow 0$ as $n \rightarrow \infty$. In particular, deduce that $E(u_n) \rightarrow 0$. Now, since all the previous equations and inequalities involving u_n hold also with u_n replaced by $u_n - E(u_n)$, deduce that $\{u_n - E(u_n)\}_n \subset \mathcal{M} \cap \mathcal{T}'$, which implies that $s \in \overline{\mathcal{M} \cap \mathcal{T}'}$. Now the proof is complete. ■

5.4 Contraction property in regression-like models with missing data

With the same notation of subsection 4.2, we shall prove that the equation

$$\begin{aligned} a_1^*(X^*) &= E \left\{ E [a_1^*(X^*) \rho(Z, \theta_0) \mid W] \frac{1 - \pi(W)}{\pi(W)} \rho'(Z, \theta_0) \mid X^* \right\} \\ &\times E^{-1} \left[\frac{1}{\pi(W)} \rho(Z, \theta_0) \rho'(Z, \theta_0) \mid X^* \right] \end{aligned} \quad (39)$$

has a unique solution which can be obtained by successive approximation, under the additional assumption

$$\inf_w \pi(w) = 1 - \beta > 0, \quad (40)$$

the infimum being taken over all possible values of W . For simplicity, in the reminder of this subsection we drop the arguments of the functions. Let $\tilde{\rho} = \pi^{-1/2} \rho$. Assuming that $E(\tilde{\rho} \tilde{\rho}' \mid X^*)$ is invertible, equation (39) can be equivalently written under the form

$$\begin{aligned} a_1^* \tilde{\rho} &= E \left[E(a_1^* \rho \mid W) \frac{1 - \pi}{\pi} \rho' \mid X^* \right] E^{-1} \left(\frac{1}{\pi} \rho \rho' \mid X^* \right) \tilde{\rho} \\ &= E [E(a_1^* \tilde{\rho} \mid W) (1 - \pi) \tilde{\rho}' \mid X^*] E^{-1} (\tilde{\rho} \tilde{\rho}' \mid X^*) \tilde{\rho} \\ &=: \tilde{T}(a_1^* \tilde{\rho}). \end{aligned}$$

We will show that the map \tilde{T} is a contraction. Before that, let us state a Cauchy-Schwarz inequality for matrix valued random variables, a version of an inequality in Lavergne (2008): let \mathbb{E} denote the conditional expectation given an arbitrary σ -field, let $A \in \mathbb{R}^n \times \mathbb{R}^p$ and $B \in \mathbb{R}^n \times \mathbb{R}^q$ be random matrices such that $\mathbb{E}(\text{tr}(A'A)), \mathbb{E}(\text{tr}(B'B)) < \infty$ and $\mathbb{E}(A'A)$ is non-singular. Then the matrix $\mathbb{E}(B'B) - \mathbb{E}(B'A)\mathbb{E}^{-1}(A'A)\mathbb{E}(A'B)$ is positive semi-definite, with equality if and only if $B = A\mathbb{E}^{-1}(A'A)\mathbb{E}(A'B)$.⁴ We also use the following notation: for any symmetric matrices B_1, B_2 , $B_1 \gg B_2$ means $B_1 - B_2$ is positive semi-definite. Let us write

$$\begin{aligned}
E[\tilde{T}(a_1^* \tilde{\rho}) \tilde{T}'(a_1^* \tilde{\rho})] &= E \left\{ [E(a_1^* \tilde{\rho} | W) (1 - \pi) \tilde{\rho}' | X^*] E^{-1}(\tilde{\rho} \tilde{\rho}' | X^*) \tilde{\rho} \right. \\
&\quad \left. \times \tilde{\rho}' E^{-1}(\tilde{\rho} \tilde{\rho}' | X^*) \{ [E(a_1^* \tilde{\rho} | W) (1 - \pi) \tilde{\rho}' | X^*] \}' \right\} \\
&= E \left\{ [E(a_1^* \tilde{\rho} | W) (1 - \pi) \tilde{\rho}' | X^*] E^{-1}(\tilde{\rho} \tilde{\rho}' | X^*) \right. \\
&\quad \left. \times \{ [E(a_1^* \tilde{\rho} | W) (1 - \pi) \tilde{\rho}' | X^*] \}' \right\} \\
(\text{Cauchy-Schwarz}) &\ll E \left\{ E \left[E(a_1^* \tilde{\rho} | W) (1 - \pi)^2 E(\tilde{\rho}' a_1^* | W) | X^* \right] \right\} \\
&= E \left[E(a_1^* \tilde{\rho} | W) (1 - \pi)^2 E(\tilde{\rho}' a_1^* | W) \right] \\
(\text{Cauchy-Schwarz}) &\ll E \left[(1 - \pi)^2 (a_1^* \tilde{\rho}) (a_1^* \tilde{\rho})' \right]
\end{aligned}$$

This implies

$$\begin{aligned}
\left\| \tilde{T}(a_1^* \tilde{\rho}) \right\|_{L^2}^2 &= E \left\{ \text{tr} \left[\tilde{T}'(a_1^* \tilde{\rho}) \tilde{T}(a_1^* \tilde{\rho}) \right] \right\} = \text{tr} \left\{ E \left[\tilde{T}(a_1^* \tilde{\rho}) \tilde{T}'(a_1^* \tilde{\rho}) \right] \right\} \\
&\leq \sup_w [1 - \pi(w)] \|a_1^* \tilde{\rho}\|_{L^2}^2 \leq \beta \|a_1^* \tilde{\rho}\|_{L^2}^2,
\end{aligned}$$

where $\beta = \sup_w [1 - \pi(w)] = 1 - \inf_w \pi(w) < 1$ by assumption (40). Deduce that \tilde{T} is a contracting map.

⁴Like in Lavergne (2008), let $\Lambda = \mathbb{E}^{-1}(A'A)\mathbb{E}(A'B)$. Then

$$\mathbb{E}[(B - A\Lambda)'(B - A\Lambda)] = \mathbb{E}(B'B) - \mathbb{E}(B'A)\mathbb{E}^{-1}(A'A)\mathbb{E}(A'B)$$

is clearly positive semi-definite, and is zero iff $B = A\Lambda$.

5.5 Efficient score with parametric selection probability in regression-like models with missing data

Let $X^{(1)} = X^*$, $X^{(2)} = W$ and the parameter vector $\theta = (\alpha', \gamma)'$. Moreover, let

$$\begin{aligned} g_1(Z, \theta) &= \frac{\delta}{\pi(W, \gamma)} \rho(Y, X^*, \alpha), & g_2(Z, \theta) &= \frac{\delta}{\pi(W, \gamma)} - 1, \\ \bar{S}_\theta &= \bar{a}_1(X^*) g_1(Z, \theta) + \bar{a}_2(W) g_2(Z, \theta) = \begin{pmatrix} \bar{S}_\alpha \\ \bar{S}_\gamma \end{pmatrix}, \end{aligned}$$

where

$$\begin{aligned} \bar{a}_1(X^*) &= \bar{a}_1(X^{(1)}) \\ &= \begin{pmatrix} -\partial_\alpha E [E(\pi^{-1}(W, \gamma_0)\delta | X^*, W) \rho' | X^*] E^{-1}(\pi^{-1}(W, \gamma_0)\rho \rho' | X^*) \\ 0 \end{pmatrix} \\ &+ E \{ E[\bar{a}_1(X^*) \rho | W] (\pi^{-1}(W, \gamma_0) - 1) \rho' | X^* \} E^{-1}(\pi^{-1}(W, \gamma_0)\rho \rho' | X^*). \end{aligned}$$

If we partition $\bar{a}_1(X^*)$ in $\bar{a}_1(X^*) = \begin{pmatrix} \bar{a}_{1,\alpha}(X^*) \\ \bar{a}_{1,\gamma}(X^*) \end{pmatrix}$ and we use the same short notation as previously, the preceding equations can be written as

$$\begin{aligned} \bar{a}_{1,\alpha}(X^*) &= -\partial_\alpha E \left(E \left(\frac{\delta}{\pi} | X^*, W \right) \rho' | X^* \right) E^{-1} \left(\frac{1}{\pi} \rho \rho' | X^* \right) \\ &+ E \left\{ E[\bar{a}_{1,\alpha}(X^*) \rho | W] \left(\frac{1}{\pi} - 1 \right) \rho' | X^* \right\} \\ &\quad \times E^{-1} \left(\frac{1}{\pi} \rho \rho' | X^* \right), \\ \bar{a}_{1,\gamma}(X^*) &= E \left\{ E[\bar{a}_{1,\gamma}(X^*) \rho | W] \left(\frac{1}{\pi} - 1 \right) \rho' | X^* \right\} \\ &\quad \times E^{-1} \left(\frac{1}{\pi} \rho \rho' | X^* \right), \end{aligned}$$

with the obvious solution $\bar{a}_{1,\gamma} \equiv 0$ for the subvector of \bar{a}_1 corresponding to γ (possibly not the unique solution, but any solution yields the same efficient score \bar{S}_θ). Similar calculations can be done for $\bar{a}_2(W)$:

$$\begin{aligned}
\bar{a}_2(W) &= \bar{a}_2(X^{(2)}) \\
&= \begin{pmatrix} 0 \\ \frac{1}{\pi} \partial_\gamma \pi \end{pmatrix} \frac{\pi}{1-\pi} - E[\bar{a}_1(X^*) \rho | W],
\end{aligned}$$

which gives, for $\bar{a}_2(W) = \begin{pmatrix} \bar{a}_{2,\alpha}(W) \\ \bar{a}_{2,\gamma}(W) \end{pmatrix}$,

$$\bar{a}_{2,\alpha}(W) = -E[\bar{a}_{1,\alpha}(X^*) \rho | W]$$

$$\bar{a}_{2,\gamma}(W) = \frac{1}{1-\pi} \partial_\gamma \pi - E[\bar{a}_{1,\gamma}(X^*) \rho | W] = \frac{1}{1-\pi} \partial_\gamma \pi.$$

Therefore,

$$\bar{S}_\theta = \bar{a}_1(X^*) g_1 + \bar{a}_2(W) g_2 = \begin{pmatrix} \bar{S}_\alpha \\ \bar{S}_\gamma \end{pmatrix} = \begin{pmatrix} \bar{a}_{1,\alpha}(X^*) g_1 + \bar{a}_{2,\alpha}(W) g_2 \\ \bar{a}_{2,\gamma}(W) g_2 \end{pmatrix},$$

where

$$\begin{aligned}
\bar{a}_{1,\alpha}(X^*) &= -\partial_\alpha E(\rho' | X^*) E^{-1} \left(\frac{1}{\pi} \rho \rho' | X^* \right) \\
&\quad + E \left\{ E[\bar{a}_{1,\alpha}(X^*) \rho | W] \frac{1-\pi}{\pi} \rho' | X^* \right\} \\
&\quad \times E^{-1} \left(\frac{1}{\pi} \rho \rho' | X^* \right),
\end{aligned}$$

$$\bar{a}_{2,\alpha}(W) = -E[\bar{a}_{1,\alpha}(X^*) \rho | W],$$

$$\bar{a}_{2,\gamma}(W) = \frac{\pi}{1-\pi} \partial_\gamma \pi \quad \left(= \frac{\pi(W, \gamma_0)}{1-\pi(W, \gamma_0)} \partial_\gamma \pi(W, \gamma_0) \right).$$

Now, for any $s = b(W) \cdot g_2 = b(W) \left(\frac{\delta}{\pi(W, \gamma_0)} - 1 \right) \in \mathcal{T}_2^\perp$, we have

$$\begin{aligned}
E(\bar{S}_\alpha s' | W) &= E \left[\bar{S}_\alpha \left(\frac{\delta}{\pi(W, \gamma_0)} - 1 \right) | W \right] b'(W) \\
&= E \left\{ \left[\bar{a}_{1,\alpha}(X^*) \frac{\delta}{\pi} \rho + \bar{a}_{2,\alpha}(W) \left(\frac{\delta}{\pi} - 1 \right) \right] \left(\frac{\delta}{\pi} - 1 \right) | W \right\} b'(W) \\
&= \{ E[\bar{a}_{1,\alpha}(X^*) \rho | W] + \bar{a}_{2,\alpha}(W) \} \left(\frac{1}{\pi} - 1 \right) b'(W) \\
&= \{ E[\bar{a}_{1,\alpha}(X^*) \rho | W] - E[\bar{a}_{1,\alpha}(X^*) \rho | W] \} \left(\frac{1}{\pi} - 1 \right) b'(W) \\
&= 0,
\end{aligned}$$

so that, since $\bar{S}_\gamma = \bar{a}_{2,\gamma}(W) \cdot g_2$, we obtain

$$E(\bar{S}_\alpha \bar{S}_\gamma') = E \left[E(\bar{S}_\alpha \bar{S}_\gamma' | W) \right] = 0.$$

This means that the efficient score S_α^* for α , equal to the residual of the (componentwise) projection of \bar{S}_α on \bar{S}_γ , coincides with \bar{S}_α ,

$$S_\alpha^* = \bar{S}_\alpha - E(\bar{S}_\alpha \bar{S}_\gamma') V^{-1}(\bar{S}_\gamma) \bar{S}_\gamma = \bar{S}_\alpha,$$

and has the same expression, as already noticed in Robins, Rotnitzky and Zhao (1994), as in the case where $\pi(W)$ is completely known :

$$\begin{aligned}
S_\alpha^* &= \bar{S}_\alpha = \bar{a}_{1,\alpha}(X^*) g_1 + \bar{a}_{2,\alpha}(W) g_2 \\
&= a_1^*(X^*) g_1 + a_2^*(W) g_2.
\end{aligned}$$