



Statistical Theory and Related Fields

ISSN: 2475-4269 (Print) 2475-4277 (Online) Journal homepage: https://www.tandfonline.com/loi/tstf20

### An equivalence result for moment equations when data are missing at random

Marian Hristache & Valentin Patilea

To cite this article: Marian Hristache & Valentin Patilea (2019) An equivalence result for moment equations when data are missing at random, Statistical Theory and Related Fields, 3:2, 199-207, DOI: 10.1080/24754269.2019.1672021

To link to this article: https://doi.org/10.1080/24754269.2019.1672021

4	1	(	1
Е			
Е			

Published online: 09 Oct 2019.



Submit your article to this journal 🗗

Article views: 27



View related articles



View Crossmark data 🗹

Citing articles: 1 View citing articles

### An equivalence result for moment equations when data are missing at random

Marian Hristache<sup>†</sup> and Valentin Patilea<sup>†</sup>

Univ Rennes, Ensai, CNRS, CREST-UMR 9194, Rennes, France

#### ABSTRACT

We consider general statistical models defined by moment equations when data are missing at random. Using the inverse probability weighting, such a model is shown to be equivalent with a model for the observed variables only, augmented by a moment condition defined by the missing mechanism. Our framework covers a large class of parametric and semiparametric models where we allow for missing responses, missing covariates and any combination of them. The equivalence result is stated under minimal technical conditions and sheds new light on various aspects of interest in the missing data literature, as for instance the efficiency bounds and the construction of the efficient estimators, the restricted estimators and the imputation.

1. Introduction

Models defined by moment and conditional moment equations are widely used in statistics, biostatistics and econometrics; see, for instance, Ai and Chen (2003, 2012), Domínguez and Lobato (2004), and the references therein. Here, we investigate general moment or conditional moment equation models with missing data. The main idea we propose is that under a missing at random assumption, the initial model with missing data is equivalent with a inverse probability weighting moment equations model for the complete observations, augmented by a moment condition defined by the missing mechanism. The equivalence, a generalisation of the GMM equivalence result of Graham (2011), is stated in terms of sets of probability measures. It has numerous implications and provides valuable insight, for instance on the efficiency bound calculations and the construction of efficient estimators.

In the framework of missing data, the assumption of missing at random (MAR) is presumably the most used when trying to describe an ignorable mechanism on the missingness. However, this concept, first introduced by Rubin (1976), does not have the same meaning for everyone. For simplicity, let the full observations be i.i.d. replications of a vector L = (X, Y, Z)and let  $R = (R_X, R_Y, R_Z) \in \{0, 1\}^3$  be a random vector such that its component takes the value 1 if we observe the corresponding component of L and 0 otherwise. For Rubin (1976) (see also, for example, Little & Rubin, 2002; Robins & Gill, 1997), MAR means that missingness depends only on the observed **ARTICLE HISTORY** 

Received 19 December 2018 Revised 20 September 2019 Accepted 21 September 2019

#### **KEYWORDS**

Efficiency bounds; imputation; inverse probability weighting; semiparametric regression; restricted estimators

components, denoted by  $L_{(R)}$ , of *L*:

the conditional law 
$$\mathcal{L}(R \mid L)$$
 of  $R$  given  $L$   
is the same as the conditional law  $\mathcal{L}(R \mid L_{(R)})$  of  $R$  given  $L_{(R)}$ . (1)

This concept was generalised to CAR, coarsening at random, by Heitjan and Rubin (1991) (see also, for example, van der Laan and Robins (2003)):  $\mathcal{L}(C | L)$  is the same as  $\mathcal{L}(C | \varphi(C, L))$  for an always observable transformation  $\varphi(C, L)$  of the full data L and the censoring variable C. In the context of regression-like models, the MAR assumption is usually stated in a different and more restrictive way. A strongly ignorable selection mechanism (also called conditional independence, or selection on observables, etc.) means that, assuming some components of L are always observed,

the conditional law  $\mathcal{L}(R \mid L)$  of R given L is the same

as the conditional law of R given the always observed

This assumption was originally introduced by Rosenbaum and Rubin (1983) in the framework of randomised clinical trials, which corresponds in our simple example, with L = (X, Y, Z), to the case where, for example, X is always observed, and *one and only one* of Y and Z is observed. This means that the selection vector R takes the form R = (1, D, 1 - D), where Y is observed iff D = 1 and Z is observed iff D = 0. In this

**CONTACT** Valentin Patilea valentin.patilea@ensai.fr Univ Rennes, Ensai, CNRS, CREST–UMR9194, F–35000 Rennes, France <sup>†</sup>Current address: Center for Research in Economics and Statistics (CREST), Ecole Nationale de la Statistique et de l'Analyse de l'Information (Ensai), Campus de Ker-Lann, rue Blaise Pascal, BP 37203, 35172 Bruz cedex, France.

situation, MAR means

$$P(D = 1 | X, Y, Z) = P(D = 1 | X, Y)$$
  
= 1 - P(D = 0 | X, Y, Z)  
= 1 - P(D = 0 | X, Z)  
= P(D = 1 | X, Z),

or, equivalently,

$$D \perp \!\!\!\perp Z \mid X, Y \text{ and } D \perp \!\!\!\perp Y \mid X, Z.$$
 (3)

Meanwhile a strongly ignorable missingness mechanism writes

$$P(D = 1 \mid X, Y, Z) = P(D = 1 \mid X),$$

or, equivalently,

$$D \perp (Y, Z) \mid X. \tag{4}$$

Clearly, condition (4) implies condition (3), but the reverse is not true in general. In the present work we consider the case of i.i.d. replications of a vector containing missing components for which the same subvector is missing for the incomplete replicates. In this case the MAR assumption (1) and the the strongly ignorable MAR assumption (2) coincide (and are equivalent to CAR), as is it is also the case, for example, in Cheng (1994), Tsiatis (2007), Graham (2011), among others.

Other MAR-related assumptions appear in the literature. For instance, when the response Y is missing, while X and Z are observed, Wei, Ma, and Carroll (2012) consider the assumption  $R_Y \perp (X, Y) \mid Z$  that is stronger than the MAR assumption (2), commonly used for regression models. Another assumption for the missingness mechanism is introduced in Wooldridge (2007) : W = (X, Y) and  $S \in \{0, 1\}$  is a random variable such that W and Z are observed whenever S = 1, and  $S \perp W \mid Z$ . Wooldridge's assumption is more general than the MAR condition (2) where Z is supposed to be always observed. Indeed, Wooldridge (2007) does not suppose that W and/or Z are missing if S = 0.

The paper is organised as follows. The main equivalence result is stated in Section 2. In Section 3, we revisit some examples considered in the literature in the MAR setup: estimating mean functionals in parametric and nonparametric regressions; and quantile regression with missing responses and/or covariates. For these examples, our equivalence result suggests new ways for calculating efficiency bounds and constructing efficient estimators, using for instance the GMM, empirical likelihood approaches, the SMD approach of Ai and Chen (2007), or the kernel-based method of Lavergne and Patilea (2013). In Section 4 we reinterpret some classes of so-called restricted estimators; see, for instance, Tsiatis (2007) and Tan (2011). Finally, in Section 5 we use our general result to discuss on a common belief that the (multiple) imputation is necessary in order to capture all the information from the partially observed data.

### 2. Equivalent moment model

The following statement is a version of Theorems 1 and 2 in Hristache and Patilea (2017). The proof is very similar and hence will be omitted. In the following, vectors a columns matrices and for any matrix A, A' denotes its transpose.

**Theorem 2.1:** Let  $\mathcal{M}_1$  and  $\mathcal{M}_2$  be two models defined for random vectors  $(D, W', V', U')' \in \{0, 1\} \times \mathbb{R}^{d_W} \times \mathbb{R}^{d_V} \times \mathbb{R}^{d_U}$  as follows:

$$\mathcal{M}_{1}:\begin{cases} E[\rho_{j}(\gamma, W, V, U)] = 0, \quad \forall j \in J, \\ D \perp \{U, V\} \mid W, \end{cases}$$
(5)

and

$$\mathcal{M}_{2}: \begin{cases} E\left[\frac{D}{\pi(W)}\rho_{j}(\gamma, W, V, U)\right] = 0, \quad \forall j \in J, \\ E\left[\frac{D}{\pi(W)} - 1 \mid V, W\right] = 0, \end{cases}$$

$$(6)$$

where  $\gamma \in \Gamma$  is an unknown (possibly infinite dimensional) parameter,  $\rho_j : \Gamma \times \mathbb{R}^{d_W} \times \mathbb{R}^{d_U} \times \mathbb{R}^{d_U} \to \mathbb{R}$ , for  $j \in J$ , is a collection of known measurable functions, and  $\pi$  is a unknown measurable function such that  $\pi(W) > 0$  almost surely.

The models  $M_1$  and  $M_2$  are equivalent if restricted to the laws of (D, W', V', DU')'; more precisely,

- (1)  $(D, W', V', U')' \in \mathcal{M}_1 \Rightarrow (D, W', V', U')' \in \mathcal{M}_2,$
- (2)  $(D, W', V', U')' \in \mathcal{M}_2 \Rightarrow \exists (\widetilde{D}, \widetilde{W}', \widetilde{V}', \widetilde{U}')' \in \mathcal{M}_1$  such that  $(\widetilde{D}, \widetilde{W}', \widetilde{V}', \widetilde{D}\widetilde{U}')'$  and (D, W', V', DU')' have the same distribution.

### **Remarks:**

- (1) The parameter  $\gamma$  in model  $\mathcal{M}_1$  could include parameters of interest and parameters of nuisance.
- (2) The function π(·) usually called the propensity score, could be considered completely unknown and modelled nonparametrically, or modelled using a parametric model. With at hand an estimate of π(·) obtained from the second equation in the model M<sub>2</sub>, one could use existing moment equation approaches for the estimation of the parameters in the first equation of M<sub>2</sub>. See our Example 3.1.
- (3) The link of this theorem with models where data are missing at random is made if we consider that the vector U is observed if and only if D = 1. The theorem then basically says that at the observational level, which means for the laws of the

observed vector (D, W', V', DU')', the two models  $\mathcal{M}_1$  and  $\mathcal{M}_2$  are equivalent. As a consequence, inference for the law of (D, W', V', U')' in the model  $\mathcal{M}_1$ , a moment conditions model under an assumption of data missing at random, could be done based on the model  $\mathcal{M}_2$ , which is defined using only the observed part (D, W', V', DU')' of the vector vector (D, W', V', U')'. In particular, efficiency bound calculations and efficient estimator constructions could be done in the model  $\mathcal{M}_2$ , which in many cases could be much easier.

(4) The underlying condition '*DU* is always observed' includes the usual case

D = 0 if *U* is not observed, D = 1 if *U* is observed,

but it is more general. When D = 1 one observes the value of U. Meanwhile, one should read that when D = 0, U could be observed or not since whatever the value of U is, DU = 0.

### 3. Some examples revisited

In this section we present two examples of models already studied in the literature for which our approach gives new insights and sometimes allows for simpler methods for obtaining efficiency bounds and asymptotically efficient estimators. The guiding principle is to use Theorem 2.1 and put the model of interest, in the presence of a MAR mechanism, under an equivalent form

$$E[g_1(\theta, \alpha, X, Y, Z) | X] = 0$$
  

$$E[g_2(\alpha, X, Y, Z) | X, Z] = 0,$$
(7)

where the two sets of equations are orthogonal, meaning that

$$E[g_1(\theta, \alpha, X, Y, Z)g'_2(\alpha, X, Y, Z) | X, Z] = 0.$$

The equivalent model (7) has a sequential moment structure that allows to compute the efficiency bound; see Ai and Chen (2012). Moreover, the finite dimensional interest parameter  $\theta$  can be efficiently estimated from the first equations, with the (possibly infinite dimensional) nuisance parameter  $\alpha$  known or suitably estimated from the last equations. A similar statement on the efficient estimation of  $\theta$ , in the particular case of a finite dimensional  $\alpha$  and without conditioning on *X* and *X*, *Z*, can be found in Theorem 2.2, point 8, of Prokhorov and Schmidt (2009).

## 3.1. Mean functionals with data missing at random

Consider the problem of estimating the mean of functionals of the variables in a parametric regression model with missing responses:

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

$$E[Y - r(X, \alpha) | X] = 0.$$
(8)

The parameter of interest here is  $\theta = E[h(X, Y)]$ , where  $h(\cdot, \cdot)$  is some given squared-integrable function; see Müller (2009). Hristache and Patilea (2017) considered the same framework and focused on the case where h(X, Y) does not depend on X. Here we investigate the general case where h(X, Y) that could also depend on X. Some usual examples are the mean of the response variable (h(x, y) = y), the second-order moment of the response  $(h(x, y) = \operatorname{vec}(yy'))$ , the cross-product of the response and the covariate vector  $(h(x, y) = \operatorname{vec}(yx'))$ . (Here,  $\operatorname{vec}(\cdot)$  is the vectorisation operator that transforms a matrix in a column vector by stacking the columns of the matrix.) For simplicity, we take Y with real values in the following of this section.

The regression function  $r(x, \alpha)$  has a known (parametric) form, *X* is always observed, *Y* is only observed when D = 1 and a MAR assumption holds :  $D \perp Y \mid X$ . With  $\pi(x) = P(D = 1 \mid X = x)$ , the model can be written, at the observational level, under the following equivalent form:

$$E\left\{\frac{D}{\pi(X)}\left[h(X,Y)-\theta\right]\right\} = 0$$

$$E\{D[Y-r(X,\alpha)] \mid X\} = 0$$

$$E\left[\frac{D}{\pi(X)}-1 \mid X\right] = 0.$$
(9)

The last two equations being orthogonal, since

$$E\left\{\left[\frac{D}{\pi(X)} - 1\right]D[Y - r(X,\alpha)] \mid X\right\}$$
$$= \left[\frac{1}{\pi(X)} - 1\right]E\{D[Y - r(X,\alpha)] \mid X\} = 0,$$

it is also equivalent to the model defined by the following system of orthogonal equations, where  $\sigma^2(X)$ stands for the conditional variance V(Y|X):

$$E\left\{\frac{D}{\pi(X)}\left[h(X,Y)-\theta\right]\right.$$
$$\left.-\frac{1}{\sigma^{2}(X)\pi(X)}E\left[\frac{D}{\pi(X)}h(X,Y)(Y-r(X,\alpha))\mid X\right]\right.$$
$$D[Y-r(X,\alpha)]$$
$$\left.-E\left[\frac{D}{\pi(X)}(h(X,Y)-\theta)\mid X\right]c\left[\frac{D}{\pi(X)}-1\right]\right\}=0$$
$$E\{D[Y-r(X,\alpha)]\mid X\}=0$$
$$E\left[\frac{D}{\pi(X)}-1\mid X\right]=0.$$
(10)

Solving for  $\theta$ , we get

$$\theta = E[\Phi(Y, X, D; \alpha, \sigma^2, \pi, \eta_1, \eta_2)]$$

where

$$\Phi(Y, X, D; \alpha, \sigma^{2}, \pi, \eta_{1}, \eta_{2})$$

$$= \frac{D}{\pi(X)}h(X, Y) - E\left[\frac{D}{\pi(X)}h(X, Y) \mid X\right]$$

$$\times \left[\frac{D}{\pi(X)} - 1\right]$$

$$- \frac{1}{\sigma^{2}(X)\pi(X)}E\left[\frac{D}{\pi(X)}h(X, Y)(Y - r(X, \alpha)) \mid X\right]$$

$$\times D[Y - r(X, \alpha)],$$

$$\eta_{1}(X) = E[Dh(X, Y) \mid X]$$

and

$$\eta_2(X) = \eta_2(X;\alpha) = E[Dh(X,Y)(Y - r(X,\alpha)) | X].$$

Let  $\hat{\alpha}$  be an estimator of  $\alpha$  obtained in the model. With the variance  $\sigma^2(\cdot)$  and the functions  $\eta_1(\cdot)$  and  $\eta_2(\cdot; \cdot)$  estimated nonparametrically, the plug-in estimator

$$\widehat{\theta} = \frac{1}{n} \sum_{i=1}^{n} \Phi(Y_i, X_i, D_i; \widehat{\alpha}, \widehat{\sigma^2}, \widehat{\pi}, \widehat{\eta_1}, \widehat{\eta_2})$$

would be efficient. Since the first equation in system (10) is orthogonalised with respect to the last one, for the propensity score  $\pi(\cdot)$ , one could use a parametric model without affecting the efficiency bound.

# 3.2. Quantile regression with data missing at random

A particular setting of quantile regression with missing data at random is considered in Wei et al. (2012). For  $0 < \tau < 1$ , the conditional quantile  $Q_{\tau}(Y | X, Z)$  of the always observed response *Y* given the regressor vectors *Z* (always observed) and *X* (observed iff D = 1) is assumed to be linear,

$$Q_{\tau}(Y | X, Z) = X' \beta_{1,\tau} + Z' \beta_{2,\tau}, \qquad (11)$$

and the missingness mechanism is defined by the strong missing at random condition

$$D \perp\!\!\!\perp (X, Y) \mid Z. \tag{12}$$

Taking in (6) U = X, V = Y, W = Z,  $\rho_j(\beta_\tau, W, V, U)$ =  $(X', Z')' [\mathbb{1}_{\{Y-X'\beta_{1,\tau}-Z'\beta_{2,\tau}\leq 0\}} - \tau] \times a_j(U, W) \triangleq$  $\rho(X, Y, Z, \beta_\tau) \times a_j(X, Z), j \in \mathbb{N}$ , where the family of functions  $\{a_j\}_{j\in\mathbb{N}}$  spans  $L^2(X, Z)$ , the model defined by (11) and (12) can be written under the following equivalent form:

$$E[D\rho(Y, X, Z, \beta_{\tau}) | X, Z] = 0$$

$$E\left[\frac{D}{\pi(Z)} - 1 | Z\right] = 0.$$
(13)

The two sets of equations being already orthogonal (with respect to the  $\sigma$ -field  $\sigma(X, Z)$ ), in this situation we can efficiently estimate the parameter  $\beta_{\tau} =$ 

 $(\beta'_{1,\tau}, \beta'_{2,\tau})'$  from the complete data only, that is from the model defined by (11) keeping for the statistical analysis only the observations for which all the components of the vector (Y, X', Z')' are observed. The gain in efficiency observed in the simulation experiment of Wei et al. (2012) for their multiple imputation improved estimator comes, in our opinion, from the supplementary parametric assumption on the form of the conditional density of *X* given *Z* (see their Assumption 4).

A more general linear quantile regression model defined by (11) with missing data at random is considered in Chen, Wan, and Zhou (2014). With their notations, we have

$$Y = Z'\theta(\tau) + \varepsilon, \quad P(\varepsilon \le 0 \mid Z) = \tau, \quad 0 < \tau < 1,$$
(14)

for the full data model. They also denote by *X* the always observed components of the vector (Y, Z')' and with  $X^c$  the components of the same vector that are observed iff the binary variable *D* takes the value 1 and use the 'standard' missing at random assumption  $P(D = 1 | Y, Z) = P(D = 1 | X, X^c) = P(D = 1 | X) = \pi(X)$ . This fits our framework by taking U = X, V = 1,  $W = X^c$  and

$$\rho_j(\theta(\tau), W, V, U) = Z[\mathbb{1}_{\{Y - Z'\theta(\tau) \le 0\}} - \tau] \times a_j(U, W)$$
$$\triangleq \rho(Y, Z, \theta(\tau)) \times a_j(Z), \quad j \in \mathbb{N},$$

where the family of functions  $\{a_j\}_{j\in\mathbb{N}}$  spans  $L^2(Z)$ . The equivalent moment equations model, at the observational level, can be written as

$$E\left\{\frac{D}{\pi(X)}Z[\mathbb{1}_{\{Y-Z'\theta(\tau)\leq 0\}}-\tau] \mid Z\right\} = 0$$

$$E\left[\frac{D}{\pi(X)}-1 \mid X\right] = 0.$$
(15)

The information bound for this model is given in Hristache and Patilea (2016). It can not be calculated explicitly, except some special cases, which includes the missing responses as before or the case where X or/and Z are discrete. It is different from the information bound given in Chen, Hong, and Tarozzi (2008) which corresponds to a model defined by an *uncon-ditional* quantile moment and a MAR assumption and could be represented equivalently under the form

$$E\left\{\frac{D}{\pi(X)}Z[\mathbb{1}_{\{Y-Z'\theta(\tau)\leq 0\}}-\tau]\right\} = 0$$

$$E\left[\frac{D}{\pi(X)}-1 \mid X\right] = 0.$$
(16)

Models (15) and (16) are quite different and so are the corresponding efficiency bounds, so that no estimation procedure given in Chen et al. (2014) could be efficient in their linear quantile regression model (14) with missing data at random.

### 4. Restricted estimators for quantile regressions and general conditional moment models with data missing at random

The model defined by the regression-like equation

$$E[\rho(\theta, Y, X, V) | X, V] = 0,$$

and the MAR selection mechanism

$$P(D = 1 | Y, X, V, W) = P(D = 1 | W) = \pi(W)$$

is equivalent, at the observational level, to the following model defined by conditional moment equations :

$$\mathcal{P}: \begin{cases} E\left[\frac{D}{\pi(W)}\rho(\theta, Y, X, V) \mid X, V\right] = 0, \\ E\left[\frac{D}{\pi(W)} - 1 \mid W\right] = 0. \end{cases}$$

This framework includes many situations. For instance, taking W' = (Y', V', Z') we obtain the case in which some regressors (conditioning variables) X are missing, while with W' = (X', V', Z') we cover the case of missing responses. Splitting Y in an observed subvector  $Y_o$  and a not always observed subvector  $Y_u$ , with  $W' = (Y'_o, V', Z')$  this corresponds to the case where both some responses and some covariates are missing. In all these examples, U is the vector of not always observed components of the data vector.

For the model

$$\mathcal{P}_{(1)}: \quad E\left[\frac{D}{\pi(W)}\rho(\theta, Y, X, V) \,|\, X, V\right] = 0,$$

denoting by  $P_0$  the true law of (Y', X', V', Z')', the tangent space is

$$\mathcal{T}_{(1)} = \left\{ s \in \{L^2(P_0)\}^{\oplus d} : E(s) = 0, \\ E\left[\frac{D}{\pi(W)}\rho(\theta, Y, X, V)s'(Y, X, V, Z) \mid X, V\right] \\ = 0 \right\}.$$

For the model

$$\mathcal{P}_{(2)}: \quad E\left[\frac{D}{\pi(W)}-1 \mid W\right]=0,$$

the tangent space is

$$\mathcal{T}_{(2)} = \left\{ s \in \{L^2(P_0)\}^{\oplus d} : E(s) = 0, \\ E\left[\left(\frac{D}{\pi(W)} - 1\right)s'(Y, X, V, Z) \mid W\right] = 0 \right\}.$$

The tangent space  $\mathcal{T}$  of  $\mathcal{P} = \mathcal{P}_{(1)} \cap \mathcal{P}_{(2)}$  is (see Hristache & Patilea, 2016)

$$\mathcal{T} = \mathcal{T}_{(1)} \cap \mathcal{T}_{(2)}.$$

We obtain the efficient score  $\overline{S}_{\theta}$  by projecting the score  $S_{\theta}$  on  $\mathcal{T}^{\perp}$ ,

$$\overline{S}_{\theta} = \Pi(S_{\theta} \mid \mathcal{T}^{\perp}) = \Pi(S_{\theta} \mid \overline{\mathcal{T}_{(1)}^{\perp} + \mathcal{T}_{(2)}^{\perp}}),$$

which gives the following solution :

$$\bar{S}_{\theta} = a_1^*(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) + a_2^*(W)$$
$$\times \left(\frac{D}{\pi(W)} - 1\right) \in \mathcal{T}_{(1)}^{\perp} + \mathcal{T}_{(2)}^{\perp},$$

where

$$\begin{aligned} a_1^*(X,V) &= \left\{ -E(\partial_\theta \rho' | X, V) \right. \\ &+ E\left[ E(a_1^*\rho \mid W) \frac{1-\pi}{\pi} \rho' \mid X, V \right] \right\} \\ &\times E^{-1}\left( \frac{1}{\pi(W)} \rho \rho' \mid X, V \right), \\ a_2^*(W) &= -E[a_1^*(X,V)\rho \mid W]. \end{aligned}$$

**Remark:**  $\overline{S}_{\theta}$  is also the efficient score in the model

$$\mathcal{P}: \begin{cases} E\left[a_1^*(X,V)\frac{D}{\pi(W)}\rho(\theta,Y,X,V)\right] = 0\\ E\left[a_2^*(W)\left(\frac{D}{\pi(W)} - 1\right)\right] = 0, \end{cases}$$

,

or in the model defined by the moment condition

$$E\left[a_1^*(X, V)\frac{D}{\pi(W)}\rho(\theta, Y, X, V) + a_2^*(W) \times \left(\frac{D}{\pi(W)} - 1\right)\right] = 0.$$

As shown in Hristache and Patilea (2016),  $a_1^*$  satisfies an equation of the form

$$a_1^*(X, V) = \gamma(X, V) + T(a_1^*(X, V)),$$

with *T* a contraction operator. The solution of this equation is unique, but in order to obtain it one needs to use nonparametric estimators at each step of the iterative procedure. An alternative approach would be to consider finite dimensional subspaces  $S_1 \subset T_{(1)}^{\perp}$  and  $S_2 \subset T_{(2)}^{\perp}$  when calculating the 'efficient score', leading to an approximately efficient score. We obtain in this way what is known in the literature as *restricted estimators*. We can write:

$$\mathcal{T}_{(1)}^{\perp} = \left\{ s = a_1(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) : a_1 \in L^2(P_0) \right\}$$

 $S_1 \subset T_{(1)}^{\perp}$  finite dimensional  $\Rightarrow \exists a_1^{(1)}, \ldots, a_1^{(k)} \in L^2(P_0)$  s.t.

$$S_{1} = \lim \left\{ a_{1}^{(i)}(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) : \\ 1 \le i \le k \right\}$$
$$\Leftrightarrow S_{1}^{\perp} = \left\{ s \in \{L^{2}(P_{0})\}^{\oplus d} : E\left(a_{1}^{(i)} \frac{D}{\pi} \rho s'\right) = 0 \\ 1 \le i \le k \right\}.$$

Compare to

$$\mathcal{T}_{(1)} = \left\{ s \in \{L^2(P_0)\}^{\oplus d} : E\left(\frac{D}{\pi}\rho s' \mid X, V\right) = 0 \right\}.$$

Similarly for  $\mathcal{S}_2 \subset \mathcal{T}_{(2)}^{\perp}$ :

$$\mathcal{T}_{(2)}^{\perp} = \left\{ s = a_2(W) \left( \frac{D}{\pi(W)} - 1 \right) : a_2 \in L^2(P_0) \right\}$$

 $S_2 \subset T_{(2)}^{\perp}$  finite dimensional  $\Rightarrow \exists a_2^{(1)}, \dots, a_2^{(l)} \in L^2(P_0)$  s.t.

$$S_2 = \ln \left\{ a_2^{(j)}(W) \left( \frac{D}{\pi(W)} - 1 \right) : 1 \le j \le l \right\}$$
  
$$\Leftrightarrow \quad S_2^{\perp} = \left\{ s \in \{L_0^2(P_0)\}^{\oplus d} : E \left[ a_2^{(j)} \left( \frac{D}{\pi(W)} - 1 \right) s' \right] = 0, \quad 1 \le j \le k \right\}.$$

An optimal class 1 restricted estimator (see Tan, 2011; Tsiatis, 2007) is solution of the approximated efficient score equation

$$E\left\{\overline{a}_{1}^{(1)}(X,V)\frac{D}{\pi(W)}\rho(\theta,Y,X,V) + \overline{a}_{2}^{(1)}(W) \times \left(\frac{D}{\pi(W)} - 1\right)\right\} = 0,$$

where  $\overline{a}_1^{(1)}$  and  $\overline{a}_2^{(2)}$  are given by

$$\overline{S}_{\theta} = \Pi(S_{\theta} | S_1 + S_2)$$
  
=  $\overline{a}_1^{(1)}(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) + \overline{a}_2^{(1)}(W)$   
 $\times \left(\frac{D}{\pi(W)} - 1\right).$ 

In fact,  $\overline{S}_{\theta}$  is the efficient score in the following moment equations model:

$$\mathcal{P}': \begin{cases} E\left[a_{1}^{(1)}(X,V)\frac{D}{\pi(W)}\rho(\theta,Y,X,V)\right] = 0\\ \vdots\\ E\left[a_{1}^{(k)}(X,V)\frac{D}{\pi(W)}\rho(\theta,Y,X,V)\right] = 0\\ E\left[a_{2}^{(1)}(W)\left(\frac{D}{\pi(W)} - 1\right)\right] = 0\\ \vdots\\ E\left[a_{2}^{(l)}(W)\left(\frac{D}{\pi(W)} - 1\right)\right] = 0 \end{cases}$$

This allows for a new, simple and intuitive interpretation of the optimal class 1 restricted estimators as efficient estimators in a larger model, obtained from the initial one by using appropriate 'instruments' to transform the conditional moment equations in a (growing) number of unconditional moment conditions. Another advantage of this new perspective is the access to the most commonly used methods of obtaining efficient estimators in moment equations models such as GMM, SMD (see Lavergne & Patilea, 2013) or empirical likelihood estimators.

Similar procedures can be used for class 2 restricted estimators, based on

$$\Pi(S_{\theta} \mid S_{1} + T_{(2)}^{\perp}) = \overline{a}_{1}^{(2)}(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) + \overline{a}_{2}^{(2)}(W) \left(\frac{D}{\pi(W)} - 1\right)$$

and class 3 restricted estimators (Tan, 2011), based on

$$\Pi(S_{\theta} \mid \mathcal{T}_{(1)}^{\perp} + S_2) = \overline{a}_1^{(3)}(X, V) \frac{D}{\pi(W)} \rho(\theta, Y, X, V) + \overline{a}_2^{(3)}(W) \left(\frac{D}{\pi(W)} - 1\right).$$

### 4.1. Simulation study

The approach on restricted estimators is illustrated in a setting already considered by Chen, Wan, and Zhou (2015); see their *Example 1*, scenario  $S_2$ . With the notations of the previous section, the data are generated from the following model:

$$Y = \theta_0 + \theta_1 X + \theta_2 V + 0.5[1 + (X + V)]\varepsilon,$$
  

$$\varepsilon \sim \mathcal{N}(0, 1), \qquad (17)$$

where  $(\theta_0, \theta_1, \theta_2) = (1, -1, 1)$  and (X, V) follows a centred bivariate normal distribution with unit variances and correlation equal to 0.5. The parameter of interest here is the vector coefficient  $(\theta_0(\tau), \theta_1(\tau), \theta_2(\tau))$  of the

conditional quantile of *Y* given *X* and *V*:

$$Q_{\tau}(Y \mid X, V) = \theta_0(\tau) + \theta_1(\tau)X + \theta_2(\tau)V, \quad (18)$$

with  $\theta_0(\tau) = \theta_0 + Q_\tau(\varepsilon)$ ,  $\theta_1(\tau) = \theta_1 + 0.5Q_\tau(\varepsilon)$ ,  $\theta_2(\tau) = \theta_2 + 0.5Q_\tau(\varepsilon)$ , where  $Q_\tau(\varepsilon)$  is the  $\tau$  th quantile of  $\varepsilon$ ,  $\tau \in (0, 1)$ . Herein, we only report the case  $\tau = 0.75$ . The variables *Y* and *V* are always observed, while *X* is observed if and only if D = 1, where *D* is a Bernoulli random variable such that  $P(D = 1 \mid Y, V) = 0.4(1 + \sin^2(Y - V))\mathbb{1}_{\{|Y - V| \le 1\}} + 1 - \mathbb{1}_{\{|Y - V| \le 1\}} = \pi(W)$ , with W = (Y, V). The model for the fully observed data is defined by the regression-like equation

$$E[\rho(\theta, Y, X, V) | X, V] = 0,$$

where  $\rho(\theta, Y, X, V) = \mathbb{1}_{\{Y-\theta_0-\theta_1X-\theta_2V \le 0\}} - \tau$ . Under the MAR selection mechanism

$$P(D = 1 | Y, X, V) = P(D = 1 | Y, V) = \pi(W).$$

it is equivalent, at the observational level, to the following model defined by conditional moment equations:

$$\mathcal{P}: \begin{cases} E\left[\frac{D}{\pi(W)}(\mathbb{1}_{\{Y=\theta_0-\theta_1X=\theta_2V\leq 0\}}-\tau)\mid X, V\right]=0,\\ E\left[\frac{D}{\pi(W)}-1\mid W\right]=0. \end{cases}$$

The restricted estimators considered are obtained by the generalised method of moments in the following models  $\mathcal{P}_s$ ,  $s \in \{a, b, c, d, e, f\}$ , which contain the model  $\mathcal{P}$ :

$$\mathcal{P}_{s}: \begin{cases} E\left[\frac{D}{\pi_{s}(\phi, Y, V)}(\mathbb{1}_{\{Y-\theta_{0}-\theta_{1}X-\theta_{2}V\leq 0\}}-\tau) \\ a_{1s}^{(k)}(X, V)\right] = 0, \quad k \in \{1, \dots, k_{s}\} \\ E\left\{\left[\frac{D}{\pi_{s}(\phi, Y, V)}-1\right]a_{2s}^{(l)}(Y, V)\right\} \\ = 0, \quad l \in \{1, \dots, l_{s}\}, \end{cases}$$

where:

- (a)  $\pi_a \equiv 1, D \equiv 1, k_a = 3, a_{1a}^{(1)}(X, V) = 1, a_{1a}^{(2)}(X, V)$ =  $X, a_{1a}^{(3)}(X, V) = V$  (no missing data, 1, X and V as instruments);
- (b)  $\pi_b \equiv 1, D \equiv 1, k_b = 3, a_{1b}^{(1)}(X, V) = (1 + X^2)^{-1}, a_{1b}^{(2)}(X, V) = (1 + V^2)^{-1}, a_{1b}^{(3)}(X, V) = (1 + |X| + |V|)^{-2}$  (no missing data,  $(1 + X^2)^{-1}, (1 + V^2)^{-1}$ and  $(1 + |X| + |V|)^{-2}$  as instruments);
- (c)  $\pi_c(Y, V) = 0.4(1 + \sin^2(Y V))\mathbb{1}_{\{|Y-V| \le 1\}} + 1 \mathbb{1}_{\{|Y-V| \le 1\}}, k_c = 3, a_{1c}^{(1)}(X, V) = (1 + X^2)^{-1}, a_{1c}^{(2)}(X, V) = (1 + V^2)^{-1}, a_{1c}^{(3)}(X, V) = (1 + |X| + |V|)^{-2}$  (true propensity score,  $(1 + X^2)^{-1}, (1 + V^2)^{-1}$  and  $(1 + |X| + |V|)^{-2}$  as instruments);
- (d)  $\pi_d(Y, V) = \{1 + \exp[-(\phi_0 + \phi_1 Y + \phi_2 V)]\}^{-1},$ with  $\phi_0$ ,  $\phi_1$  and  $\phi_2$  estimated from a logistic regression,  $k_d = 3$ ,  $a_{1d}^{(1)}(X, V) = (1 + X^2)^{-1},$

 $a_{1d}^{(2)}(X, V) = (1 + V^2)^{-1}, a_{1d}^{(3)}(X, V) = (1 + |X| + |V|)^{-2}$  (propensity score estimated by a logistic regression on *Y* and *V*,  $(1 + X^2)^{-1}$ ,  $(1 + V^2)^{-1}$  and  $(1 + |X| + |V|)^{-2}$  as instruments for the IPW quantile equation);

- (e)  $\pi_e(Y, V) = \{1 + \exp[-(\phi_0 + \phi_1 Y + \phi_2 V)]\}^{-1}, k_e = 3, a_{1e}^{(1)}(X, V) = (1 + X^2)^{-1}, a_{1e}^{(2)}(X, V) = (1 + V^2)^{-1}, a_{1e}^{(3)}(X, V) = (1 + |X| + |V|)^{-2}, l_e = 3, a_{2e}^{(1)}(Y, V) = 1, a_{2e}^{(2)}(Y, V) = Y, a_{2e}^{(3)}(Y, V) = V, ((1 + X^2)^{-1}, (1 + V^2)^{-1} \text{ and } (1 + |X| + |V|)^{-2} \text{ as instruments for the IPW quantile equation, 1, Y and V as instruments for the propensity score equation);$
- (f)  $\pi_f(Y, V) = \{1 + \exp\{-[\phi_0 + \phi_1(Y V) + \phi_2 (Y V)^2]\}\}^{-1}, k_f = 3, a_{1f}^{(1)}(X, V) = (1 + X^2)^{-1}, a_{1f}^{(2)}(X, V) = (1 + V^2)^{-1}, a_{1f}^{(3)}(X, V) = (1 + |X| + |V|)^{-2}, l_f = 3, a_{2f}^{(1)}(Y, V) = 1, a_{2f}^{(2)}(Y, V) = Y V, a_{2f}^{(3)}(Y, V) = (Y V)^2, ((1 + X^2)^{-1}, (1 + V^2)^{-1} and (1 + |X| + |V|)^{-2} as instruments for the IPW quantile equation, 1, Y V and (Y V)^2 as instruments for the propensity score equation).$

The estimates of the MSE obtained in the case  $\tau = 0.75$  from 1000 replications, with sample size  $n \in \{200, 400, \dots, 1400, 1600\}$ , are given in Table 1.

Note that none of the GMM estimators in models  $\mathcal{P}_s$  could be efficient in the initial model  $\mathcal{P}$ , but only approximately efficient, if the instruments are suitably chosen, which could be a delicate point in practice. Here we observe that the instruments  $a_{1b}^{(k)}$ , involved in the first equations in the first equations of the model  $\mathcal{P}_b$  performs better that the instruments  $a_{1a}^{(k)}$  used in  $\mathcal{P}_a$ . We observe a similar phenomenon for the propensity score equations when looking at the columns  $\mathcal{P}_d$ ,  $\mathcal{P}_e$  and  $\mathcal{P}_f$ . The case in  $\mathcal{P}_d$  corresponds to common practice when one trusts the logistic regression for the propensity score. The cases in  $\mathcal{P}_e$  and  $\mathcal{P}_f$  correspond to our approach based on instruments with more effective instruments in the later case. The non-orthogonality of the quantile model equations and the propensity score equations could explain the better results in  $\mathcal{P}_f$ . A joint estimation of the two set of equations with effective

**Table 1.** Estimates of  $E(\|\widehat{\theta} - \theta\|^2)$  over 1000 replicates when  $\tau = 0.75$ .

n	$\mathcal{P}_{a}$	$\mathcal{P}_b$	$\mathcal{P}_{c}$	$\mathcal{P}_d$	$\mathcal{P}_{e}$	$\mathcal{P}_{f}$
200	431.29	221.19	233.81	404.44	333.66	246.55
400	350.79	127.95	143.11	240.37	275.73	118.66
600	296.33	101.49	113.57	213.83	241.23	85.69
800	276.33	82.23	81.28	196.20	232.01	73.81
1000	257.09	73.58	73.77	176.90	211.63	65.14
1200	240.27	66.14	66.49	177.00	208.15	62.44
1400	237.85	59.56	66.47	173.35	196.65	58.42
1600	231.86	53.85	56.53	163.77	192.95	52.06

Note: The reported values are multiplied by 10<sup>4</sup>.

instruments could improve over the common practice. Next, let us notice that the models  $\mathcal{P}_c$  and  $\mathcal{P}_f$  are similar: we use the same instruments for the first equation in  $\mathcal{P}_s$ . Moreover, in  $\mathcal{P}_c$  we use the true propensity score, while in  $\mathcal{P}_f$  we use an estimated propensity score obtained from a model that is somehow close to the true propensity score. As the two equations in the model  $\mathcal{P}_s$  are not orthogonal, estimating the propensity score could improve the asymptotic variance of the estimators of  $\theta$ . This is related to the so-called puzzling phenomenon noticed by Prokhorov and Schmidt (2009). Here, even if propensity score the model is slightly wrong, there is still a gain of MSE. Let us also note the surprisingly good results for the model  $\mathcal{P}_f$  in which  $\log[\pi/(1 - \pi)]$  $\pi$ )] is approximated by a quadratic function of *Y* –*V*. Using the same instruments for the conditional quantile equations, the estimation with missing data is even better than in the case with full data (compare results for model  $\mathcal{P}_f$  to those for model  $\mathcal{P}_b$ ). This could be explained by the fact that in model  $\mathcal{P}_b$  we do not use the optimal instruments that should be proportional to the conditional density of the error term at the origin. The weighting introduced by the propensity score seems, in some sense, to compensate the non-optimal instruments. This suggests further possible improvements based on other choices of instrumental variables in order to approach efficiency.

### 5. Is imputation really informative?

Multiple imputation is a widely used method to generate substitute values when data are missing. However, under the MAR assumption, the interest of multiple imputation in the context of conditional moment restriction models is at least questionable, as discussed in the following.

Consider that (D, W', V', DU')' is always observed and consider the MAR assumption

$$(U,V) \perp D \mid W. \tag{19}$$

Then, any substitute observation generated from the law of  $\widetilde{U}$  is adequate to replace a missing U, where the law of  $\widetilde{U}$  should be such that

$$\mathcal{L}(\widetilde{U} \mid \widetilde{W}, \widetilde{V}, \widetilde{D} = 0) = \mathcal{L}(U \mid W, V, D = 1)$$
$$= \mathcal{L}(\widetilde{U} \mid \widetilde{W}, \widetilde{V}, \widetilde{D} = 1).$$

(Here,  $\mathcal{L}(V_1 \mid V_2)$  denotes the conditional law of  $V_1$  given  $V_2$ .) Since, in general, the law  $\mathcal{L}(U \mid W, V, D = 1)$  is unknown, one can estimate it, parametrically or nonparametrically, and generate substitute observations from this estimate. This is the so-called parametric or nonparametric imputation. See, for instance, Wang and Chen (2009), Wei et al. (2012), Chen and Van Keilegom (2013) for some nonparametric imputation applications.

The equivalence established by Theorem 2.1 for models defined by moment restrictions, implies that *all* the information on the parameter  $\theta$  in the initial model under the MAR assumption (19) is contained in the model defined by the equations (6). Let us point out that the last equation of the model (6) includes the information contained in the incomplete observations. Indeed, to estimate  $\pi(\cdot)$ , parametrically or nonparametrically, one uses *all* the observations of *W*. This remark opens new perspectives for defining estimators of  $\theta$  without using substitute observations. Moreover, this remark sheds some new light on a common justification used in the literature, namely that imputation is necessary in order to capture the information contained in the partially observed data.

### 6. Conclusions

We consider a statistical model defined by an arbitrary number of moment equations. Our framework includes a large panel of models defined through conditional and/or unconditional moments. Next, we assume that some variables are missing at random. In this setup of modelling with missing data, we present a model equivalence result. It states that the initial statistical model together with the MAR mechanism is equivalent to a moment equations model. Using the equivalent model could greatly simplify the estimation and the inference with missing data problems. We discuss several consequences for widely used models, including the quantile regressions.

### **Disclosure statement**

No potential conflict of interest was reported by the authors.

### Notes on contributors

*Marian Hristache* is Associated Professor, Ecole Nationale de la Statistique et de l'Analyse de l'Information (Ensai), Rennes, France (E-mail: marian.hristache@ensai.fr).

*Valentin Patilea* is Professor, Ecole Nationale de la Statistique et de l'Analyse de l'Information (Ensai), and Center for Research in Economics and Statistics (CREST), Rennes, France (E-mail: valentin.patilea@ensai.fr).

### References

- Ai, C., & Chen, X. (2003). Efficient estimation of models with conditional moment restrictions containing unknown functions. *Econometrica*, 71, 1795–1843.
- Ai, C., & Chen, X. (2007). Estimation of possibly misspecified semiparametric conditional moment restriction models with different conditioning variables. *Journal of Econometrics*, 141, 5–43.
- Ai, C., & Chen, X. (2012). The semiparametric efficiency bound for models of sequential moment restrictions containing unknown functions. *Journal of Econometrics*, 170, 442–457. Thirtieth Anniversary of Generalized Method of Moments.

- Chen, X., Hong, H., & Tarozzi, A. (2008). Semiparametric efficiency in GMM models with auxiliary data. *The Annals of Statistics*, *36*, 808–843.
- Chen, S. X., & Van Keilegom, I. (2013). Estimation in semiparametric models with missing data. *Annals of the Institute of Statistical Mathematics*, 65, 785–805.
- Chen, X., Wan, A. T. K., & Zhou, Y. (2014). Efficient quantile regression analysis with missing observations. *Journal of the American Statistical Association*, 110(510), 723–741.
- Chen, X., Wan, A. T. K., & Zhou, Y. (2015). Efficient quantile regression analysis with missing observations. *Journal of the American Statistical Association*, 110, 723–741.
- Cheng, P. E. (1994). Nonparametric estimation of mean functionals with data missing at random. *Journal of the American Statistical Association*, 89, 81–87.
- Domínguez, M. A., & Lobato, I. N. (2004). Consistent estimation of models defined by conditional moment restrictions. *Econometrica*, 72, 1601–1615.
- Graham, B. S. (2011). Efficiency bounds for missing data models with semiparametric restrictions. *Econometrica*, 79, 437–452.
- Heitjan, D. F., & Rubin, D. B. (1991). Ignorability and coarse data. *The Annals of Statistics*, *19*, 2244–2253.
- Hristache, M., & Patilea, V. (2016). Semiparametric efficiency bounds for conditional moment restriction models with different conditioning variables. *Econometric Theory*, 32, 917–946.
- Hristache, M., & Patilea, V. (2017). Conditional moment models with data missing at random. *Biometrika*, 104, 735–742.
- 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.

- Little, R., & Rubin, D. (2002). *Statistical analysis with missing data*. Wiley series in probability and mathematical statistics. Probability and mathematical statistics. John Wiley & Sons, Inc., Hoboken, New Jersey.
- Müller, U. U. (2009). Estimating linear functionals in nonlinear regression with responses missing at random. *The Annals of Statistics*, 37, 2245–2277.
- Prokhorov, A., & Schmidt, P. (2009). GMM redundancy results for general missing data problems. *Journal of Econometrics*, 151, 47–55.
- Robins, J. M., & Gill, R. D. (1997). Non-response models for the analysis of non-monotone ignorable missing data. *Statistics in Medicine*, 16, 39–56.
- Rosenbaum, P. R., & Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, *70*, 41–55.
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63, 581–592.
- Tan, Z. (2011). Efficient restricted estimators for conditional mean models with missing data. *Biometrika*, 98, 663–684.
- Tsiatis, A. (2007). Semiparametric theory and missing data. New York: Springer-Verlag.
- van der Laan, M. J., & Robins, J. M. (2003). Unified methods for censored longitudinal data and causality. New York: Springer-Verlag.
- Wang, D., & Chen, S. X. (2009). Empirical likelihood for estimating equations with missing values. *The Annals of Statistics*, 37, 490–517.
- Wei, Y., Ma, Y., & Carroll, R. J. (2012). Multiple imputation in quantile regression. *Biometrika*, 99, 423–438.
- Wooldridge, J. M. (2007). Inverse probability weighted estimation for general missing data problems. *Journal of Econometrics*, 141, 1281–1301.